

On the elusiveness of clusters

Steven Kelk^{1*}, Celine Scornavacca², Leo van Iersel^{3**}

¹ Department of Knowledge Engineering (DKE), Maastricht University,
P.O. Box 616, 6200 MD Maastricht, The Netherlands. steven.kelk@maastrichtuniversity.nl

² Center for Bioinformatics (ZBIT), Tübingen University, Sand 14,
72076 Tübingen, Germany. scornava@informatik.uni-tuebingen.de

³ University of Canterbury, Department of Mathematics and Statistics,
Private Bag 4800, Christchurch, New Zealand. l.j.v.iersel@gmail.com

Abstract. Rooted phylogenetic networks are often used to represent conflicting phylogenetic signals. Given a set of clusters, a network is said to represent these clusters in the *softwired* sense if, for each cluster in the input set, at least one tree embedded in the network contains that cluster. Motivated by parsimony we might wish to construct such a network using as few reticulations as possible, or minimizing the *level* of the network, i.e. the maximum number of reticulations used in any “tangled” region of the network. Although these are NP-hard problems, here we prove that, for every fixed $k \geq 0$, it is polynomial-time solvable to construct a phylogenetic network with level equal to k representing a cluster set, or to determine that no such network exists. However, this algorithm does not lend itself to a practical implementation. We also prove that the comparatively efficient CASS algorithm correctly solves this problem (and also minimizes the reticulation number) when input clusters are obtained from two not necessarily binary gene trees on the same set of taxa but does not always minimize level for general cluster sets. Finally, we describe a new algorithm which generates in polynomial-time all binary phylogenetic networks with exactly r reticulations representing a set of input clusters (for every fixed $r \geq 0$).

1 Introduction

The traditional abstraction for modeling evolution is the phylogenetic tree. The underlying principle of such a tree is that the observed diversity in a set of species (or, more abstractly, a set of *taxa*) can be explained by branching events that cause lineages to split into two or more sublineages [21,6,7]. However, there is increasing attention for the situation when observed data cannot satisfactorily be modeled by a tree. The field of phylogenetic networks has arisen with this challenge in mind. Phylogenetic networks generalize phylogenetic trees, but within this very general characterization there are many different definitions and models [12]. In this article we are concerned with *rooted* phylogenetic networks. Such networks assume that the observed data evolves from a unique starting point (the root) and that evolution is directed away from this root. The main way these networks differ from rooted phylogenetic trees is the presence of *reticulation* nodes: nodes with indegree 2 or higher. For the remainder of this article we will use the term phylogenetic network, or just network, to refer to rooted phylogenetic networks. We refer the reader to [12,19,20,11,25] for detailed background information.

Constructing a phylogenetic network that “explains” the observed data is a trivial problem if no optimality criteria are imposed upon the constructed network. One simple optimality criterion that has attracted a great deal of attention in the literature, *reticulation minimization*, is to compute a phylogenetic network that explains the observed data but using as few reticulation events (essentially, reticulation nodes) as possible. This is an algorithmically hard problem, irrespective of the exact construction technique that is being applied [25]. A related optimality criterion, *level minimization* [16,15,24,26,22] is motivated by the observation that a phylogenetic network can be

* Steven Kelk was partly funded by a Computational Life Sciences grant of The Netherlands Organisation for Scientific Research (NWO)

** Leo van Iersel was funded by the Allan Wilson Centre for Molecular Ecology and Evolution.

regarded as some kind of tree backbone decorated with *tangles* of reticulate activity [12]. Here the challenge is to construct a phylogenetic network that explains the observed data but such that the maximum number of reticulation events inside any biconnected component, the *level*, is as low as possible. Reticulation minimization is thus a global optimality criterion, and level minimization is in some sense a local optimality criterion; see Figure 1. Both criteria will have an important role in this article.

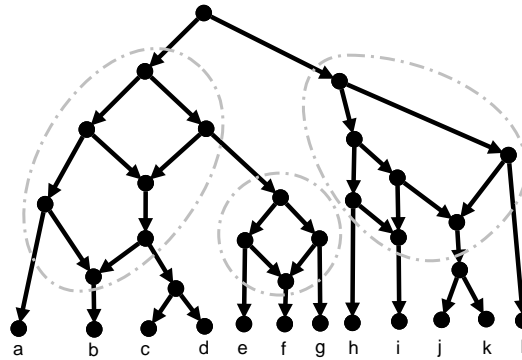


Fig. 1. Example of a phylogenetic network with five reticulations. The encircled subgraphs form its biconnected components, also known as its “tangles”. This binary network has level equal to 2 since each biconnected component contains at most two reticulations.

The question remains, when does a phylogenetic network explain the observed data? This depends very much on the exact construction technique being applied. The classical problem is motivated by the biological observation that, although the evolution of a set of organisms might best be explained by a phylogenetic network, the individual genes of the organisms will generally undergo treelike evolution [19]. In such a case one can think of the gene trees as being *displayed by* (i.e. topologically embedded within) the species network. The reticulation nodes then have an explicit biological interpretation as (for example) hybridization, recombination or horizontal gene transfer events. Hence the following problem: given a set of rooted phylogenetic trees, all on the same set of taxa, compute a phylogenetic network with a minimum number of reticulations that displays all the input trees. The problem is already NP-hard (and APX-hard) for the case of two input trees [3]. However, extensive research by different authors has shown that, by exploiting the fixed parameter tractability of the problem [1,2], the two-tree problem can be solved to satisfaction for many instances [4,31]. The case of more than two input trees, or input trees that are not all binary (i.e. some nodes have outdegree three or higher), has been considerably less well studied [29,14].

A parallel, and related, line of research concerns “piecewise” assembly of phylogenetic networks. Whereas the tree problem described above concerns the combination of a small number of large hypotheses (e.g. gene trees) into a phylogenetic network, an alternative strategy is to combine a large number of small hypotheses into a phylogenetic network. Examples of such small hypotheses include *rooted triplets* (phylogenetic trees defined on size-3 subsets of the taxa) [28,26], (binary) characters (e.g. whether or not the taxon is vertebrate) [8,9,30,18], and *clusters* (clades) [10,11,27]. Proponents of such piecewise assembly techniques argue that in this way it is easier (than with trees) to discard parts of the input that are not well-supported. In this article we focus specifically on clusters, although the classical tree problem and all the other piecewise construction techniques do play a secondary role. This secondary role is linked to the fact, as observed in [25], that under certain circumstances all these different models behave in a unified way. We shall return to this point later.

Let us then say more about the cluster model. A cluster C is a subset of the taxa and we say that a phylogenetic network *represents* the cluster in the softwired sense if *some* tree embedded

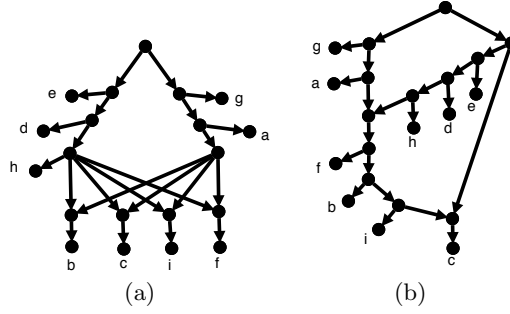


Fig. 2. (a) The output of the galled network algorithm [11] for $\mathcal{C} = \{\{a, b, f, g, i\}, \{a, b, c, f, g, i\}, \{a, b, f, i\}, \{b, c, f, i\}, \{c, d, e, h\}, \{d, e, h\}, \{b, c, f, h, i\}, \{b, c, d, f, h, i\}, \{b, c, i\}, \{a, g\}, \{b, i\}, \{c, i\}, \{d, h\}\}$ and (b) a (simple) network with two fewer reticulations that also represents this set of clusters.

in the network contains a clade equal to that cluster [11]. In other words: *some* tree T embedded in the network has an edge such that C is exactly the set of all taxa reachable from the head of that edge by directed paths. The general problem is, given a set of clusters, to construct an optimal phylogenetic network that represents all the input clusters. The set of input clusters can be constructed in an ad-hoc fashion, but often the set of clusters is generated by extracting the set of clusters induced by a set of rooted phylogenetic trees and then possibly excluding weakly supported clusters. This is the technique applied in the program DENDROSCOPE [13]. A disadvantage of this technique is that some of the topology of the original trees can be lost [25], but on the plus side it permits a focus on only well-supported clades, which is a major concern of practicing phylogeneticists.

In [11] it was shown, given a set of clusters, how to construct a *galled network* with a small number of reticulations that represents the clusters. However, given that galled networks are a restricted subclass of phylogenetic networks, it was unclear how far that algorithm actually minimizes the number of reticulations (or the level) when ranging over the entire space of phylogenetic networks, see for example Figure 2. This is the context in which the CASS algorithm was developed [27]. The CASS algorithm, in some sense a natural follow-up to the algorithm of [11], was formally shown to produce solutions of minimum *level* whenever the minimum level is at most two. However, the optimality of the CASS algorithm for “higher-level” inputs, and the performance of the algorithm in terms of minimizing number of reticulations, remained unclear. On the practical side the good news was that CASS produced solutions with fewer reticulations and lower level than the algorithm from [11]. Intriguingly it was also observed in [27] that, for several sets of input clusters induced by two binary trees, the networks produced by CASS had an identical number of reticulations to networks generated by algorithms that aim to display the trees themselves. This observation was the inspiration behind [25] in which it was proven that, in the case of two binary trees, the choice of construction technique (tree, triplets, characters, clusters) does not affect the number of reticulations required. However, this unification was shown to break down for data obtained from three or more binary trees.

After [25] several important questions about clusters remained open. Does CASS always minimize level? If not, can we find a different algorithm that efficiently minimizes level? Under which circumstances does CASS also minimize the number of reticulations? In how far do the unification results of [25] hold for non-binary trees?

The results in this article settle many of these open questions. Firstly, we show that in the case of clusters obtained from two not necessarily binary trees, a divide and conquer algorithm using CASS as subroutine, called here CASS^{DC} and implemented in DENDROSCOPE [13], does minimize both the number of reticulations, and the level. Spin-off results from this include non-binary versions of several unification results from [25], culminating in the observation that, in the case of clusters obtained from two trees, CASS^{DC} also computes the minimum number of reticulations

required to display the trees themselves (in the sense of [17]), rather than just the clusters from the trees. We also obtain deeper insights into why the two-tree case is so special and not representative for the problem on three or more trees. In particular, the two-tree case seems to be best understood as the only point at which a very natural lower bound is guaranteed to be tight.

Secondly we show that CASS does not, unfortunately, always minimize level when the input data requires solutions of level 3 or higher. We give an explicit counterexample and explain what goes wrong with the CASS algorithm in this case.

To offset this negative result we describe a polynomial-time algorithm that shows, for every fixed natural number k , how to determine whether a set of clusters can be represented by a network with level k . This algorithm, which is very different to CASS, is purely theoretical but does give important insights into the underlying structure of the cluster model. Also on the positive side we show that, for sets of clusters induced by arbitrarily large sets of *binary* trees, a simple polynomial-time algorithm can construct *all* binary phylogenetic networks with r reticulations that represent the clusters, for every fixed $r \geq 0$. To demonstrate an important design principle first observed in [25] we give a practical implementation of this algorithm, CLUSTISTIC, which elegantly “bootstraps” an existing software package for merging rooted triplets into a phylogenetic network.

To summarize, the results in this article help advance our understanding of the cluster model considerably. Nevertheless, many questions remain, and in the final section of this article we discuss a number of them. Perhaps the biggest question, which is the motivation for the title of this article, concerns the fact that the cluster model so far has not enjoyed the same kind of steady algorithmic improvements witnessed in the tree literature. Why does it seem harder to work with the clusters inside the trees than the trees themselves? To what, exactly, can the elusiveness of clusters be attributed?

2 Preliminaries

Consider a set \mathcal{X} of taxa. A *rooted phylogenetic network* (on \mathcal{X}), henceforth *network*, is a directed acyclic graph with a single node with indegree zero (the *root*), no nodes with both indegree and outdegree equal to 1, and leaves bijectively labeled by \mathcal{X} . In this article we identify the leaves with \mathcal{X} . The indegree of a node v is denoted $\delta^-(v)$ and v is called a *reticulation* if $\delta^-(v) \geq 2$. An edge (u, v) is called a *reticulation edge* if its target node v is a reticulation and is called a *tree edge* otherwise. When counting reticulations in a network, we count reticulations with more than two incoming edges more than once because, biologically, these reticulations represent several reticulate evolutionary events. Therefore, we formally define the *reticulation number* of a network $N = (V, E)$ as

$$r(N) = \sum_{v \in V: \delta^-(v) > 0} (\delta^-(v) - 1) = |E| - |V| + 1 .$$

A *rooted phylogenetic tree* on \mathcal{X} , henceforth *tree*, is simply a network that has reticulation number zero. We say that a network N on \mathcal{X} *displays* a tree T if T can be obtained from N by performing a series of node and edge deletions and eventually by suppressing nodes with both indegree and outdegree equal to 1. We assume without loss of generality that each reticulation has outdegree at least one. Consequently, each leaf has indegree one. We say that a network is *binary* if every reticulation node has indegree 2 and outdegree 1 and every tree node that is not a leaf has outdegree 2.

Proper subsets of \mathcal{X} are called *clusters*, and a cluster C is a *singleton* if $|C| = 1$. We say that an edge (u, v) of a tree *represents* a cluster $C \subset \mathcal{X}$ if C is the set of leaf descendants of v . A tree T represents a cluster C if it contains an edge that represents C . It is well-known that the set of clusters represented by a tree is a laminar set, and uniquely defines that tree. We say that a network N represents a cluster $C \subset \mathcal{X}$ “in the hardwired sense” if there exists a tree edge (u, v) of N such that C is the set of leaf descendants of v . Alternatively, we say that N represents C “in the softwired sense” if N displays some tree T on \mathcal{X} such that T represents C . In this article we only consider the softwired notion of cluster representation and henceforth assume this implicitly. A network represents a set of clusters \mathcal{C} if it represents every cluster in \mathcal{C} (and possibly more). The

set of softwired clusters of a network can be obtained as follows. For a network N , we say that a *switching* of N is obtained by, for each reticulation node, deleting all but one of its incoming edges. Given a network N and a switching T_N of N , we say that an edge (u, v) of N represents a cluster C w.r.t. T_N if (u, v) is an edge of T_N and C is the set of leaf descendants of v in T_N . The set of softwired clusters of N is the set of clusters represented by all edges of N w.r.t. T_N , where T_N ranges over all possible switchings [12]. It is also natural to define that an edge (u, v) of N represents a cluster C if there exists some switching T_N of N such that (u, v) represents C w.r.t. T_N . Note that, in general, an edge of N might represent multiple clusters, and a cluster might be represented by multiple edges of N .

Given a set of clusters \mathcal{C} on \mathcal{X} , throughout the article we assume that, for any taxon $x \in \mathcal{X}$, \mathcal{C} contains at least one cluster C containing x . For a set \mathcal{C} of clusters on \mathcal{X} we define $r(\mathcal{C})$ as $\min\{r(N) \mid N \text{ represents } \mathcal{C}\}$, we sometimes refer to this as the *reticulation number* of \mathcal{C} . The related concept of *level* requires some more background. A directed acyclic graph is *connected* (also called “weakly connected”) if there is an undirected path (ignoring edge orientations) between each pair of nodes. A node (edge) of a directed graph is called a *cut-node* (*cut-edge*) if its removal disconnects the graph. A directed graph is *biconnected* if it contains no cut-nodes. A biconnected subgraph B of a directed graph G is said to be a *biconnected component* if there is no biconnected subgraph $B' \neq B$ of G that contains B . A phylogenetic network is said to be a *level- $\leq k$ network* if each biconnected component has reticulation number less than or equal to k .⁴ A level- $\leq k$ network is called a *simple level- $\leq k$ network* if the removal of a cut-node or a cut-edge creates two or more connected components of which at most one is non-trivial (i.e. contains at least one edge). A (simple) level- $\leq k$ network N is called a (simple) *level- k network* if the maximum reticulation number among the biconnected components of N is precisely k . For example, the network in Figure 1 is a level-2 network while the one in Figure 2(a) is a simple level-4 network. Note that a tree is a level-0 network. For a set \mathcal{C} of clusters on \mathcal{X} we define $\ell(\mathcal{C})$, the *level* of \mathcal{C} , as the smallest $k \geq 0$ such that there exists a level- k network that represents \mathcal{C} . It is immediate that for every cluster set \mathcal{C} $r(\mathcal{C}) \geq \ell(\mathcal{C})$, because a level- k network always contains at least one biconnected component containing k reticulations.

We say that two clusters $C_1, C_2 \subset \mathcal{X}$ are *compatible* if either $C_1 \cap C_2 = \emptyset$ or $C_1 \subseteq C_2$ or $C_2 \subseteq C_1$. Consider a set of clusters \mathcal{C} . The *incompatibility graph* $IG(\mathcal{C})$ of \mathcal{C} is the undirected graph (V, E) that has node set $V = \mathcal{C}$ and edge set $E = \{\{C_1, C_2\} \mid C_1 \text{ and } C_2 \text{ are incompatible clusters in } \mathcal{C}\}$. We say that a set of taxa $\mathcal{X}' \subseteq \mathcal{X}$ is *separated* (by \mathcal{C}) if there exists a cluster $C \in \mathcal{C}$ that is incompatible with \mathcal{X}' , and *unseparated* otherwise.

We say that a set of clusters \mathcal{C} on \mathcal{X} is *separating* if it separates all sets of taxa \mathcal{X}' such that $\mathcal{X}' \subset \mathcal{X}$ and $|\mathcal{X}'| \geq 2$. We say that a set of clusters \mathcal{C} on \mathcal{X} is *tangled* [12] if:

1. $IG(\mathcal{C})$ is connected and has more than one node;
2. every pair of taxa x and y in \mathcal{X} is separated by \mathcal{C} .

Remember that here we assume that any taxon of \mathcal{X} is contained in at least one cluster $C \in \mathcal{C}$. Then, it can easily be verified that, given a tangled set of clusters \mathcal{C} on \mathcal{X} , \mathcal{C} is separating.

The incompatibility graph and the concept of tangled clusters are important because they highlight an important difference between (the computation of) $r(\mathcal{C})$ and $\ell(\mathcal{C})$. In [27] the authors show that, if $\ell(\mathcal{C}) = k$, then a level- k network that represents \mathcal{C} can be constructed by combining in polynomial time simple level- $\leq k$ networks constructed independently for each connected component of $IG(\mathcal{C})$. The actual procedure is slightly more involved but it shows in any case that a polynomial-time algorithm for constructing simple level- $\leq k$ networks can easily be extended to a polynomial-time algorithm for constructing level- $\leq k$ networks. We will make use of this fact in Section 3.

Unfortunately, as has been observed by several authors, the same procedure does *not* necessarily lead to networks that have reticulation number $r(\mathcal{C})$. In other words, computation of $r(\mathcal{C})$ requires something more complicated than independently optimizing each connected component of $IG(\mathcal{C})$.

⁴ Note that to determine the reticulation number of a biconnected component, the indegree of each node is computed using only edges belonging to this biconnected component.

An important special case, however, is when \mathcal{C} is separating; in this case any network N that represents \mathcal{C} is simple (or can be trivially modified to become simple) and $r(\mathcal{C}) = \ell(\mathcal{C})$. We will formalize this in due course.

To conclude the preliminaries we note that, throughout the article, we often write that an algorithm is “polynomial time” without formally specifying what the input size is. Unless otherwise specified the input is a set of clusters \mathcal{C} on taxa set \mathcal{X} . It is sufficient to take $|\mathcal{C}| + |\mathcal{X}|$ as a lower bound on the size of the input. In some cases (such as Lemma 3 in Section 3) $|\mathcal{C}|$ is at most a constant factor larger than $|\mathcal{X}|$ and then it is sufficient to prove a running time polynomial in $|\mathcal{X}|$. In other cases a running time of the form $O(|\mathcal{C}|^a |\mathcal{X}|^b)$ is obtained, for constants a and b , and this is clearly polynomial in $|\mathcal{C}| + |\mathcal{X}|$ because $(|\mathcal{C}| + |\mathcal{X}|)^2 \geq |\mathcal{C}| |\mathcal{X}|$.

2.1 Structure of the article

To facilitate the mathematical exposition we build the results of this article up in a specific order, which differs from the order presented in the introduction. We begin with Section 3: *A theoretical polynomial-time algorithm for constructing level- k networks*, where we prove that, for every fixed $k \geq 0$, the problem of determining whether a level- k network that represents \mathcal{C} exists (and if so to construct such a network) is solvable in polynomial time. This section is, compared to the rest of the article, comparatively self-contained. In Section 4: *From theory to practice: the importance of ST-sets* we describe several fundamental properties of the ST-set, a special structure that plays a central role throughout the rest of the article. In Section 5: *Clusters obtained from sets of binary trees on \mathcal{X}* we show how, given a set \mathcal{T} of binary trees on \mathcal{X} , and for each fixed $r \geq 0$, it is possible to construct in polynomial time all binary phylogenetic networks with reticulation number r that represent all the clusters in the input trees. We also describe CLUSTISTIC, which is our implementation of this algorithm built on top of already-existing software. In Section 6: *Witnesses and a natural lower bound* we further develop the theory surrounding the ST-set, explicitly relating it to the computation of reticulation number. This is used extensively in Section 7: *The optimality and non-optimality of CASS* where we give both positive and negative results for the CASS algorithm, and in the process develop a number of powerful generalizations of unification results from [25].

3 A theoretical polynomial-time algorithm for constructing level- k networks

In this section we prove that, for every fixed $k \geq 0$, the problem of determining whether a level- k network exists that represents \mathcal{C} , and if so to construct such a network, can be solved in polynomial time.

We first require some auxiliary lemmas and definitions. The proofs are rather technical so we defer them to the appendix. For a node v let $\mathcal{X}(v) \subseteq \mathcal{X}$ be the set of all taxa reachable from v by directed paths. For an edge $e = (u, v)$ we define $\mathcal{X}(e)$ to be equal to $\mathcal{X}(v)$.

Observation 1. *Let \mathcal{C} be a separating set of clusters on \mathcal{X} . Let N be any network that represents \mathcal{C} . Then each node of N has at most one leaf child and for each cut-edge (u, v) in N , $|\mathcal{X}(v)| = 1$ or $\mathcal{X}(v) = \mathcal{X}$.*

Proof. Deferred to the appendix. □

Lemma 1. *Let \mathcal{C} be a separating set of clusters on \mathcal{X} . Let N be any network that represents \mathcal{C} . Then there exists a simple network N^* with at most one leaf-child per node such that $\ell(N^*) \leq \ell(N)$.*

Proof. Deferred to the appendix. □

Observation 1 and Lemma 1 formalize the idea that any network that represents a separating set of clusters is simple or can easily be made simple by deleting certain redundant parts of it. The following lemma shows that, in terms of minimizing reticulation number or level, we can assume without loss of generality that networks are binary.

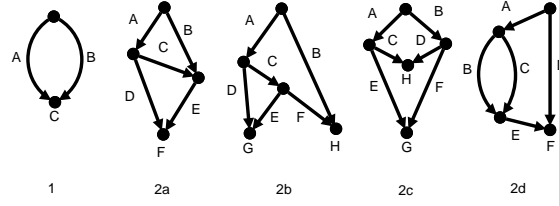


Fig. 3. The single level-1 generator and the four level-2 generators. Here the sides have been labelled with capital letters.

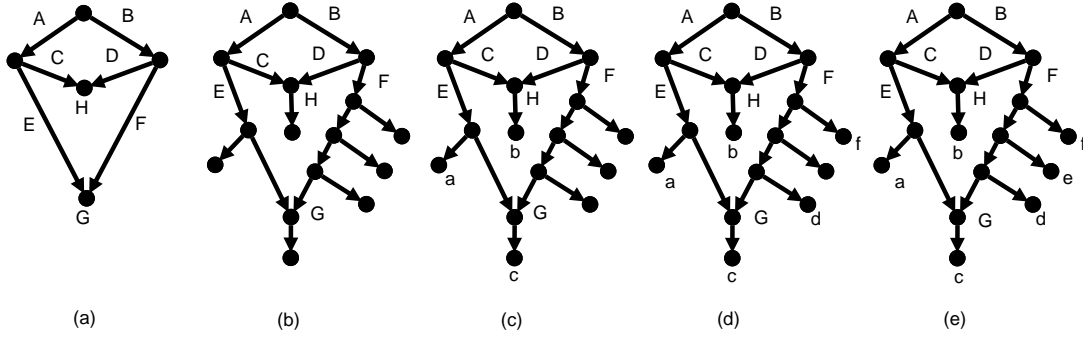


Fig. 4. An example of the execution of the algorithm outlined in Lemma 3 for the separating set of clusters $\mathcal{C} = \{\{a, b\}, \{a, c\}, \{c, d\}, \{d, e\}, \{a, b, c\}, \{c, d, e\}, \{d, e, f\}, \{c, d, e, f\}, \{b, d, e, f\}, \{b, c, d, e, f\}, \{a\}, \{b\}, \{c\}, \{d\}, \{e\}, \{f\}\}$ on $\mathcal{X} = \{a, b, c, d, e, f\}$. The value of k is fixed to 2. (a) The chosen generator g is the generator 2c in Figure 3. (b) We guess that the sides E , H and G contain one leaf while the side F contains more than 2 leaves and all other sides contain zero leaves. (c) We guess the single leaves on sides E , H and G . (d) We guess the leaves s^+ and s^- on side F . (e) We deduce the last leaf in F and we obtain a simple level-2 network on \mathcal{X} . By a stroke of luck, our first guess represents \mathcal{C} .

Lemma 2. *Let N be a phylogenetic network on \mathcal{X} . Then we can transform N into a binary phylogenetic network N' such that N' has the same reticulation number and level as N and all clusters represented by N are also represented by N' .*

Proof. Deferred to the appendix. □

Note that, given a simple binary network N , each node of N has at most one leaf child. Armed with these technical results we are ready to prove the main result of this section.

Lemma 3. *Let \mathcal{C} be separating set of clusters on \mathcal{X} . Then, for every fixed $k \geq 0$, it is possible to determine in polynomial time whether a level- k network exists that represents \mathcal{C} , and if so to construct such a network.*

Proof. From Lemmas 1 and 2 it is sufficient to focus on simple binary networks. We assume then that, for fixed k , there exists a binary simple level- k network N that represents \mathcal{C} . Let $|\mathcal{X}| = n$. Then \mathcal{C} will contain at most $2^{k+1}(n-1)$ clusters, because there are at most 2^k trees displayed by a simple level- k network, and each tree represents at most $2(n-1)$ clusters. Thus, for fixed k , the size of the input is polynomial in n . It follows from these observations that whether a set of clusters is represented by a given simple level- k network can be checked in polynomial time.

It is known that, if the leaves of N are removed and all nodes with both indegree and outdegree equal to 1 are suppressed, the resulting structure will be a level- k generator, defined in [24]. See also Figure 3. For fixed k , there are only a constant number of level- k generators [5, Proposition 2.5]. Recall that the *sides* of a level- k generator are defined as the union of its edges and its nodes

of indegree-2 and outdegree-0. For fixed k the maximum number of sides ranging over all level- k generators, is a constant.

For a cluster set \mathcal{C} on \mathcal{X} , we write $x \rightarrow y$ if and only if every non-singleton cluster in \mathcal{C} that contains x , also contains y . For example, for the cluster set of Figure 4, we have $e \rightarrow d$ and $f \rightarrow e$.

In the remainder of the proof, we illustrate a simple algorithm for determining whether a binary simple level- k network that represents \mathcal{C} exists by attempting to reconstruct such a network. Let g be the generator underlying N . We only require polynomially many tries to compute g , because there are only a constant number of generators. So assume we know g . For each side of g , we guess whether there are 0, 1, 2 or more than 2 leaves on that side. For each side containing exactly one leaf, we guess what that is. For each side s of g containing 2 or more leaves, we guess the leaf s^+ that is nearest to the root on that side, and the leaf s^- that is furthest from the root on that side. For an example see Figure 4. Note that, since for each side we have only four options (0, 1, 2 or more than 2 leaves), and in the latter case only s^+ and s^- have to be chosen, it follows that we have a polynomial number of guesses to try.

We will now show how to add the remaining leaves. Note that we may fail to insert all leaves in the network. This means that we made the wrong guess and that another set of guesses has to be checked. We say that a side s is *lowest* if it does not yet have all its leaves, and there is no other such side s' reachable from s . By reachable we mean that in the underlying generator g , there is a directed path from the head of side s to the tail of side s' . Since N is a directed acyclic graph, until all leaves in \mathcal{X} have been added, there will always be a lowest side. For example, the side F in Figure 4(d) is lowest. The idea is to add leaves to the lowest side s , until all its leaves have been added. We then continue with remaining lowest sides until we have reconstructed N .

Given a lowest side s , with s^+ and s^- fixed, it is possible to tell in polynomial time what the correct remaining leaves for s are, as follows. Observe that a leaf x that is on side s in N and which has not yet been added has the property $s^+ \rightarrow x \rightarrow s^-$. Furthermore, there is at least one cluster $C \in \mathcal{C}$ such that $\{x, s^+, s^-\} \cap C = \{x, s^-\}$. There exists at least one such cluster because otherwise $\{x, s^+\}$ would be not separated in \mathcal{C} , a contradiction since \mathcal{C} is separating. We call such a cluster a *split cluster for side s* . Now, observe that for every split cluster C for side s , and for every side $t \neq s$ that contains 2 or more leaves in N , either $\{t^+, t^-\} \cap C = \{t^+, t^-\}$ or $\{t^+, t^-\} \cap C = \emptyset$. This follows because the only edges in N that represent C lie on side s . If this is not the case, our set of guesses was incorrect and a new one has to be checked.

Now, consider any leaf y that has not yet been added to the network. Assume that this leaf belongs to side t for some t . We want to have a simple test to avoid wrongly placing it on side s , with $s \neq t$. Side t will contain three or more leaves in N , so we can assume that t^+ and t^- exist. If $s^+ \rightarrow y \rightarrow s^-$ does not hold then it is immediately clear that y cannot be put on side s . So assume (conversely) that this condition does hold, and for the same reason assume there is a split cluster C for side s that contains y . In other words, there is a cluster C such that $\{y, s^+, s^-\} \cap C = \{y, s^-\}$. Since $t^+ \rightarrow y \rightarrow t^-$ holds, it follows that C also contains t^+ and t^- , because any cluster that contains y also contains t^- , and we know that C contains either both of t^+ and t^- , or neither of them. However, there is no edge in N that can represent C : the only edges that represent C lie on side s , but the fact that s is the lowest side means that no cluster beginning on side s can contain any leaves on side t . To summarize, we have a simple test for determining whether a leaf should be placed on side s . Once we have determined the set of leaves that should be placed on side s , it is easy to determine the correct order of those leaves by inspecting \rightarrow relationships. Indeed, if q and p are two leaves that belong to side s , and q is nearer to s^- in the network we aim to reconstruct, then obviously $p \rightarrow q$. Since \mathcal{C} is separating there will be (by separation) some cluster that contains q but not p , so $q \not\rightarrow p$. If all leaves can be added in such a way, we obtain a simple level- k network on \mathcal{X} and we can check in polynomial time if it represents \mathcal{C} . If it is not the case or we fail to insert at least one leaf in the network, another set of guesses can be checked until a simple level- k network representing \mathcal{C} (if any exist) is found. This concludes the proof. \square

The following corollary follows automatically from the generality of the proof of Lemma 3.

Corollary 1. *Let \mathcal{C} be a separating set of clusters on \mathcal{X} . Then, for every fixed $k \geq 0$, it is possible to construct in polynomial time all binary simple level- k networks that represent \mathcal{C} .*

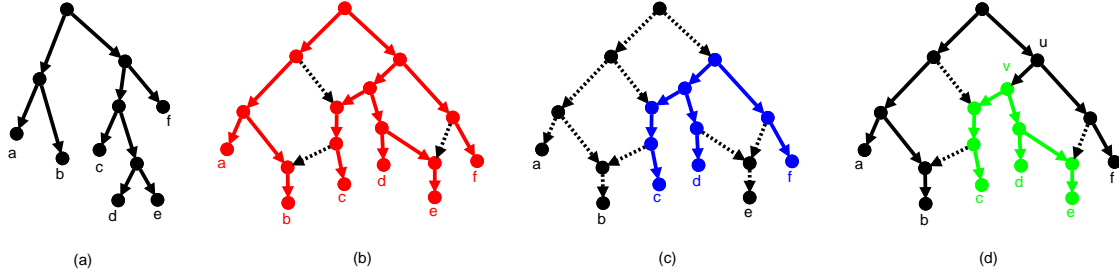


Fig. 5. A phylogenetic tree T (a) and a phylogenetic network N (b,c,d); (b) illustrates in red that N displays T (deleted edges are dashed); (c) illustrates that N is consistent with (amongst others) the triplet $cd|f$ (deleted edges are again dashed); (d) illustrates that N represents (amongst others) cluster $\{c, d, e\}$ in the softwired sense (dashed reticulation edges are “switched off”).

Using Lemma 3 we can prove the following result:

Theorem 1. *Let \mathcal{C} be a (not necessarily separating) set of clusters on \mathcal{X} . Then, for every fixed $k \geq 0$, it is possible to determine in polynomial time whether a level- k network exists that represents \mathcal{C} , and if so to construct such a network.*

Proof. Recall that all tangled cluster sets are separating. It was shown in [27] that the existence of a polynomial-time algorithm for constructing a level- $\leq k$ network from a tangled cluster set, is sufficient to give a polynomial-time algorithm for constructing level- k networks from general cluster sets. (Specifically, several tangled cluster sets are obtained by processing each non-trivial connected component of the incompatibility graph of the original cluster set [27,12]). Hence we can assume without loss of generality that \mathcal{C} is tangled. Lemma 3 is thus sufficient, and we are done. \square

3.1 Rooted triplets

It is interesting to note that the proof technique used in Lemma 3 leads to a simplified proof, presented in the following corollary, of a complexity result that was first proven in [28]. (The algorithm in [28] yielded a much faster running time, however). Let us first recall several definitions related to rooted triplets. A (*rooted*) *triplet* on \mathcal{X} is a binary phylogenetic tree on a size-3 subset of \mathcal{X} . We use $xy|z$ to denote the triplet with taxa x, y on one side of the root and z on the other side of the root. For triplets, the notion of “represent” can be formalized by the notion of “display” introduced above. However, for triplets “consistent with” is often used instead of “displayed by”. A triplet $xy|z$ is *consistent* with a phylogenetic network N (and N is *consistent* with $xy|z$) if $xy|z$ is displayed by N . See Figure 5 for an example. Given a phylogenetic tree T on \mathcal{X} , we let $Tr(T)$ denote the set of all rooted triplets on \mathcal{X} that are consistent with T . For a set of phylogenetic trees \mathcal{T} , we let $Tr(\mathcal{T})$ denote the set of all rooted triplets that are consistent with some tree in \mathcal{T} , i.e. $Tr(\mathcal{T}) = \bigcup_{T \in \mathcal{T}} Tr(T)$. A set of triplets on \mathcal{X} is *dense* if, for every size-3 subset $\{x, y, z\} \subseteq \mathcal{X}$, at least one of $xy|z$, $xz|y$, $yz|x$ is in the triplet set.

Corollary 2. *Let R be a dense set of triplets on \mathcal{X} . Then, for every fixed $k \geq 0$, it is possible to determine in polynomial time whether a binary simple level- k network exists that is consistent with R , and if so to construct such a network.*

Proof. As pointed out in [28] it is possible to determine in polynomial time whether a given network is indeed consistent with a set of input triplets. Then, the proof of Lemma 3 holds here almost entirely. The only significant difference concerns the adding of leaves to the lowest side: a not yet allocated leaf x belongs on lowest side s if and only if the triplet $s^-x|s^+$ is in the input. \square

We shall return to rooted triplets again later in the article.

4 From theory to practice: the importance of ST-sets

The algorithm described in Section 3 is polynomial time but only of theoretical interest because its running time is too high to be useful in practice. In the rest of this article we will focus on practical polynomial-time algorithms. In all these algorithms the ST-set, which can informally be thought of as treelike subsets of \mathcal{X} , has a central role. We begin by formally defining ST-sets and describing their basic properties. We will expand upon these basic properties in subsequent sections of the article.

4.1 Definition and basic properties of ST-sets

Given a set $S \subseteq \mathcal{X}$ of taxa, we use $\mathcal{C} \setminus S$ to denote the result of removing all elements of S from each cluster in \mathcal{C} and we use $\mathcal{C}|S$ to denote $\mathcal{C} \setminus (\mathcal{X} \setminus S)$ (the restriction of \mathcal{C} to S). We say that a set $S \subseteq \mathcal{X}$ is an *ST-set* with respect to \mathcal{C} , if S is not separated by \mathcal{C} and any two clusters $C_1, C_2 \in \mathcal{C}|S$ are compatible. Note that, unlike in [27], we allow the possibility that $S = \emptyset$ or $S = \mathcal{X}$. (We say that an ST-set S is *trivial* if $S = \emptyset$ or $S = \mathcal{X}$). An ST-set S is *maximal* if there is no ST-set T with $S \subset T$.

Informally, the maximal ST-sets are the result of repeatedly collapsing pairs of unseparated taxa for as long as possible; we can think of them as “islands of laminarity” within the cluster set. ST-sets first explicitly appeared in [27] but, as we shall see in due course, they implicitly arose earlier in the *recombination network* literature. An important feature of ST-sets is that there can in general be very many of them. For example, suppose \mathcal{C} contains only $|\mathcal{X}| = n$ singleton clusters; then \mathcal{C} has 2^n ST-sets. However, as the following technical results show, \mathcal{C} will have at most n *maximal* ST-sets, and they will partition \mathcal{X} i.e. they are mutually disjoint and entirely cover \mathcal{X} . Several of the proofs have been deferred to the appendix.

Lemma 4. *Let \mathcal{C} be a set of clusters on \mathcal{X} and let $S_1 \neq S_2$ be two ST-sets of \mathcal{C} . If $S_1 \cap S_2 \neq \emptyset$ then $S_1 \cup S_2$ is an ST-set.*

Proof. Deferred to the appendix. □

Corollary 3. *Let \mathcal{C} be a set of clusters on \mathcal{X} and let $S_1 \neq S_2$ be two maximal ST-sets of \mathcal{C} . Then $S_1 \cap S_2 = \emptyset$.*

Corollary 4. *Let \mathcal{C} be a set of clusters on \mathcal{X} . Then there are at most n maximal ST-sets with respect to \mathcal{C} , they are uniquely defined and they partition \mathcal{X} .*

Proof. The disjointness of maximal ST-sets guarantees that there at most n of them. Consider the set \mathcal{S} of all (necessarily disjoint) maximal ST-sets of \mathcal{X} . Suppose the maximal ST-sets in \mathcal{S} do not entirely cover \mathcal{X} . Then there is some $x \in \mathcal{X}$ which is disjoint from all maximal ST-sets in \mathcal{S} . Let S be the ST-set of largest cardinality that contains x ; such an ST-set must exist because $\{x\}$ is an ST-set. $S \notin \mathcal{S}$ so there exists some ST-set S' such that $S \subset S'$, but this contradicts the assumption on the cardinality of \mathcal{S} . Hence \mathcal{S} partitions \mathcal{X} , and by extension \mathcal{S} is unique. □

Lemma 5. *The maximal ST-sets of a set of clusters \mathcal{C} on \mathcal{X} can be computed in polynomial time.*

Proof. Deferred to the appendix. □

Corollary 4 and Lemma 5 are perhaps not so surprising, but we have nevertheless proven them rigorously to highlight the fact that *computing* maximal ST-sets is not a complexity bottleneck. In later sections we shall see that there is a link between NP-hardness and maximal ST-sets, but that the hardness lies in selecting certain maximal ST-sets with special properties, not in the computation of the maximal ST-sets *per se*.

Let \mathcal{T} be a set of trees, where each $T \in \mathcal{T}$ is a tree on \mathcal{X} . For a tree T we write $Cl(T)$ to denote the set of clusters induced by edges of T i.e. $C \in Cl(T)$ if and only if some edge of T represents C . We let $Cl(\mathcal{T}) = \cup_{T \in \mathcal{T}} Cl(T)$. Whenever we (reasonably) assume that all singleton clusters are

present in the input⁵, it is easy to see that every cluster set \mathcal{C} on \mathcal{X} can be written as $Cl(\mathcal{T})$ for some \mathcal{T} as follows. We take any proper coloring of $IG(\mathcal{C})$ (i.e. map the nodes of $IG(\mathcal{C})$ to colors such that no two adjacent nodes have the same color) and use the resulting colors to partition \mathcal{C} . Clusters that have been colored the same are all mutually compatible, so can be represented by a single tree corresponding to that color. Finally, whenever a subset of clusters pertaining to a color does not cover all elements of \mathcal{X} , the missing taxa can be attached to the root. An obvious corollary of this is that the chromatic number of $IG(\mathcal{C})$ is a lower bound on the cardinality of \mathcal{T} .

Whenever $\mathcal{C} = Cl(\mathcal{T})$ there is an important relationship between the nodes and edges of trees in \mathcal{T} , and the (maximal) ST-sets of \mathcal{C} . Let T be a (not necessarily binary) tree on \mathcal{X} . In Section 2 we defined when an edge of a tree represents a cluster. Here we extend this definition to *nodes* of trees. We say that a node v of T represents C if C is equal to the union of the clusters represented by some (not necessarily strict) subset of its outgoing edges. Note that if an edge (u, v) represents a cluster C then so does v .

Lemma 6. *Let $\mathcal{T} = \{T_1, \dots, T_m\}$ be a set of trees on \mathcal{X} . Let $\emptyset \subset \mathcal{X}' \subset \mathcal{X}$ be an unseparated set with respect to $Cl(\mathcal{T})$. Then for each $T_i \in \mathcal{T}$ there exists an edge e_i or a node v_i in T_i such that e_i or v_i represents \mathcal{X}' .*

Proof. Consider an arbitrary tree $T_i \in \mathcal{T}$. Every cluster in $Cl(T_i)$ is either disjoint from \mathcal{X}' , a superset of it, or a subset of it, otherwise \mathcal{X}' would be separated. If some edge of T_i represents \mathcal{X}' then we are done. Otherwise consider some node v furthest from the root which represents a cluster C such that $\mathcal{X}' \subset C$. Such a v must exist because if necessary we can take the root as v . C is equal to the union of the clusters represented by some subset of the edges outgoing from v . Each cluster represented by an outgoing edge of v is either disjoint from \mathcal{X}' or a subset of it, because of the assumption on the distance of v from the root. For the same reason, \mathcal{X}' intersects with at least two such outgoing edge clusters of v . But when \mathcal{X}' intersects with such a cluster it must contain it entirely, so \mathcal{X}' is equal to the union of some subset of the clusters represented by edges outgoing from v . Hence v represents \mathcal{X}' . \square

The following corollary is automatic.

Corollary 5. *Let $\mathcal{T} = \{T_1, \dots, T_m\}$ be a set of binary trees on \mathcal{X} . Let $\emptyset \subset \mathcal{X}' \subset \mathcal{X}$ be an unseparated set with respect to $Cl(\mathcal{T})$. Then for each $T_i \in \mathcal{T}$ there exists an edge e_i such that e_i represents \mathcal{X}' .*

Note that Lemma 6 and Corollary 5 hold in particular for (maximal) ST-sets, because all ST-sets are unseparated. Hence the two following straightforward extensions to ST-sets, which we will use extensively in the next section. Recall that an ST-set S is trivial if $S = \emptyset$ or $S = \mathcal{X}$.

Corollary 6. *Let $\mathcal{T} = \{T_1, \dots, T_m\}$ be a set of binary trees on \mathcal{X} . Let S be a non-trivial ST-set with respect to $Cl(\mathcal{T})$. Then for each $T_i \in \mathcal{T}$ there exists an edge e_i in T_i such that e_i represents S .*

Corollary 7. *Let $\mathcal{T} = \{T_1, \dots, T_m\}$ be a set of binary trees on \mathcal{X} . Then $Cl(\mathcal{T})$ contains at most $2(n-1)$ non-trivial ST-sets, and for every such ST-set S of $Cl(\mathcal{T})$ and every tree $T_i \in \mathcal{T}$ there exists a unique edge e_i of T_i such that $\mathcal{X}(e_i) = S$ and such that the subtree rooted at the head of e_i is the unique tree that represents exactly the cluster set $Cl(\mathcal{T})|S$.*

Proof. Given any tree $T_i \in \mathcal{T}$, $Cl(T_i)$ will contain exactly $2(n-1)$ edges and consequently $2(n-1)$ clusters. Since an ST-set is by definition unseparated, each ST-set of $Cl(\mathcal{T})$ is a subset of the cluster set $Cl(T_i)$ (for any i) and there are thus at most $2(n-1)$ ST-sets. Moreover, by definition, for any ST-set S of $Cl(\mathcal{T})$ we have that the set $Cl(\mathcal{T})|S$ is compatible. Since two non-isomorphic binary trees on the same taxa set induce at least two incompatible clusters, this concludes the proof. \square

⁵ The presence or absence of the singleton clusters in the input does not change the complexity of the problems we study because it is trivial to modify a network without raising its reticulation number or level such that it also represents all the singleton clusters.

Informally Corollary 7 states that when \mathcal{T} consists solely of binary trees, each non-trivial ST-set corresponds to some subtree that is common to all trees in \mathcal{T} .

5 Clusters obtained from sets of *binary* trees on \mathcal{X}

Let \mathcal{T} be a set of binary trees on \mathcal{X} . In this section we prove that, for fixed $r \geq 0$, it is possible to construct in polynomial time all binary phylogenetic networks with reticulation number r that represent $Cl(\mathcal{T})$. We also describe our new program CLUSTISTIC which implements this algorithm. CLUSTISTIC is in itself an important “proof of concept”: it has been rapidly prototyped by, building on insights from [25], slightly modifying existing software that was originally conceived to reconstruct binary level- k networks not from clusters but rooted triplets.

Let N be a network on \mathcal{X} and let T' be some tree on $\mathcal{X}' \subset \mathcal{X}$. We say that T' is a *Subtree Below a Reticulation* (SBR) of N if there is a reticulation node v in N such that no reticulation nodes $v' \neq v$ are reachable from v by a directed path, and that the subnetwork rooted at v (or, when v has outdegree exactly 1, the child of v) is exactly equal to T' . It is easy to show that (by virtue of its acyclicity) every network contains at least one SBR [28]. A simple though critical observation is:

Observation 2. *If T' is an SBR of a network N on \mathcal{X} that represents a cluster set \mathcal{C} , and T' has taxa set \mathcal{X}' , then \mathcal{X}' is an ST-set with respect to \mathcal{C} .*

We will make repeated use of this throughout the rest of the article. (Note however that in general an ST-set will not specify a unique SBR).

Given a network N with an SBR T we denote by $N \setminus_T$ the network obtained from N by deleting T and for as long as necessary applying the following tidying-up operations until they are no longer needed: deleting any node with outdegree zero that is not labelled by an element of \mathcal{X} ; suppressing all nodes with indegree and outdegree both equal to 1; replacing multi-edges with single edges; deleting nodes with indegree-0 and outdegree-1.

We call (S_1, S_2, \dots, S_p) ($p \geq 0$) a *(maximal) ST-set sequence* of \mathcal{C} if S_1 is a (maximal) ST-set of \mathcal{C} , S_2 is a (maximal) ST-set of $\mathcal{C} \setminus S_1$, S_3 is a (maximal) ST-set of $\mathcal{C} \setminus S_1 \setminus S_2$ and so on. Such a sequence is additionally a *tree* sequence if all the clusters in $\mathcal{C} \setminus S_1 \setminus \dots \setminus S_p$ are mutually compatible i.e. can be represented by a tree. We denote by p be the *length* of the sequence; $p = 0$ denotes the empty tree sequence.

Lemma 7. *Let N be a network that represents some cluster set \mathcal{C} on \mathcal{X} . Then there exists a sequence of SBRs that need to be removed to prune N into a tree, and this corresponds to an ST-set tree sequence of \mathcal{C} of length $r(N)$.*

Proof. If N is a tree then there is an empty tree sequence. Otherwise, N has an SBR T' with taxa set \mathcal{X}' and the network $N \setminus_{T'}$ represents the cluster set $\mathcal{C} \setminus \mathcal{X}'$. Clearly $r(N') < r(N)$. By observation 2, we let S_1 equal \mathcal{X}' and $d = r(N) - r(N')$. If $d > 1$ then we let S_2, \dots, S_d all equal \emptyset , the empty ST-set. We use these empty ST-sets to model the situation when removing the SBR would cause multiple reticulation nodes to disappear simultaneously. Now, N' also has at least one SBR, so we can iterate the whole process. We can repeat this until we obtain a tree: at this point we will have an ST-set tree sequence of length exactly $r(N)$. \square

In the following observation we make use of the fact that the operation $N \setminus_T$ is also meaningfully defined when the tree T is exactly the subtree rooted at some node of N .

Observation 3. *Let $\mathcal{T} = \{T_1, \dots, T_m\}$ be a set of binary trees on \mathcal{X} . Let S be a non-trivial ST-set of $Cl(\mathcal{T})$. Then $Cl(\mathcal{T}) \setminus S = Cl(\mathcal{T}')$ where \mathcal{T}' is a set of at most m binary trees $\{T'_1, \dots, T'_m\}$ on $\mathcal{X} \setminus S$ with $T'_i = T_i \setminus T_v$, where $e_i = (u, v)$ is the edge of T_i that represents S (which exists by Corollary 7) and T_v is the subtree rooted at v .*

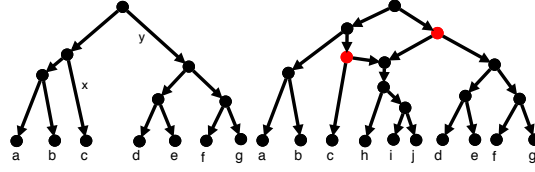


Fig. 6. Let \mathcal{C} be some set of clusters represented by the network on the right. If we guess that ST-set $S = \{h, i, j\}$ corresponds to an SBR and remove it, we obtain the tree on the left. To reverse the process we add (a tree corresponding to) S below a reticulation node whose incoming edges subdivide edges x and y .

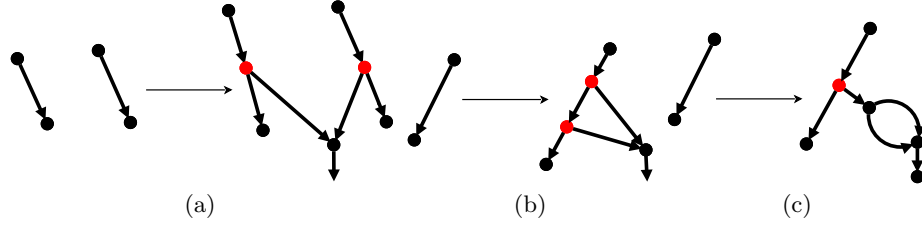


Fig. 7. The different ways of adding a reticulation back into a network, as discussed in the proof of Theorem 2. (a) Two different edges are subdivided; (b) one edge is subdivided twice; (c) one edge is subdivided once, under which a multi-edge is placed.

Theorem 2. Let $\mathcal{T} = \{T_1, \dots, T_m\}$ be a set of binary trees on \mathcal{X} . Then for a constant $r \geq 0$ it is possible to construct in polynomial time all binary networks with reticulation number at most r that represent $Cl(\mathcal{T})$ (if any exist).

Proof. Without loss of generality assume we wish to construct all such networks with reticulation number *exactly* r . Let us suppose that at least one such network N exists. By Lemma 7 there is an ST-set tree sequence $\mathcal{S} = (S_1, \dots, S_r)$ for $Cl(\mathcal{T})$ and this corresponds to a sequence of SBRs that, when removed, will prune N into a tree. Let $|\mathcal{X}| = n$. Now, note that by Observation 3 and Corollary 7 there are at most $O(n^r)$ ST-set tree sequences, which is polynomial in n for fixed constant r . \mathcal{S} will be one of these, so we can find \mathcal{S} in polynomial time. It remains to show how, assuming we have found \mathcal{S} , we can reconstruct N . First, note that for a binary tree T on \mathcal{X} , T is the *unique* tree on \mathcal{X} that represents $Cl(T)$. The clusters that still remain after removing S_r can be represented by a tree, and in particular (by Observation 3) by a unique binary tree. We call this tree N_r . We want to obtain N_{r-1} by inserting (a tree corresponding to) S_r into N_r . In particular, we wish to introduce a new reticulation node into N_r below which a tree T (itself binary and unique by Corollary 7) that represents S_r will be attached. We have two possibilities to do this: we subdivide two (not necessarily distinct) edges and use these as the tails of the new reticulation edges, see Figures 7(a) and (b), or we subdivide one edge *once* and use two *identical* reticulation edges, see Figure 7(c). Note that, if we only subdivide one edge once, we actually create a multi-edge and by our definition of phylogenetic network such edges are not allowed. However it might be necessary to create such multi-edges during *intermediate* iterations to ensure that all phylogenetic networks, including “redundant” ones, are constructed.

Unfortunately we do not know in general which edge(s) of N_r to subdivide to create the new connection point(s). However, there are only polynomially many edges in N_r , so we simply try them all. Then we repeat the process, inserting S_{r-1} into N_{r-1} to obtain N_{r-2} , and so on until we have obtained N_0 . Given that r is a constant we can test in polynomial time whether N_0 represents all the clusters in \mathcal{C} (see proof of Lemma 3).

Just as in the similar algorithm described in [28] there are several slight technicalities that should be noted. Whenever some $S_i = \emptyset$ we use a “dummy” taxon (i.e. some taxon not in \mathcal{X}) as the tree that we attach below a reticulation. The function of this is to ensure that subsequent

iterations can subdivide the edge leaving the reticulation i.e. it is a placeholder. (This will be necessary when the removal of a single SBR caused the disappearance of two or more reticulations). These dummy taxa can be removed just before N_0 is inspected to check whether it represents \mathcal{C} . Any N_0 that at this point still contains a dummy taxon whose parent is a reticulation, should be rejected, because it means at least one reticulation was not used. Secondly, when we construct the tree N_r we actually add a “dummy root” which is simply a new node ρ' and a single edge from (ρ', ρ) where ρ is the root of N_r . This deals with the situation when the removal of some SBR caused the current root to disappear and a new root to take its place. At the end, ρ' and the edge leaving it should be removed. Finally, note that any multi-edges created by intermediate iterations of the algorithm should all have disappeared (i.e. have been subdivided by reticulation edges) by the time N_0 has been reached; for this reason we reject any N_0 that still contains multi-edges.

The network N we would like to reconstruct will eventually be found as some N_0 , and given that we made no assumptions about N this shows that the algorithm constructs all possible N . \square

Corollary 8. *Let $\mathcal{T} = \{T_1, \dots, T_m\}$ be a set of binary trees on \mathcal{X} . Then for a constant $k \geq 0$ it is possible to construct in polynomial time all binary simple level- $\leq k$ networks that represent $Cl(\mathcal{T})$ (if any exist).*

Proof. For each network produced by the algorithm described in Theorem 2 we can easily check in polynomial time whether it is biconnected. \square

It is worth noting at this stage an important link with the rooted triplet literature. Recall the following proposition and lemma from [25], an article in which the relationship between trees, clusters and triplets was discussed more broadly. Proposition 1 refers to not necessarily binary trees.

Proposition 1. (Van Iersel, Kelk [25]) *For any set \mathcal{T} of trees on the same set \mathcal{X} of taxa, any phylogenetic network on \mathcal{X} representing $Cl(\mathcal{T})$ is consistent with $Tr(\mathcal{T})$.*

Lemma 8. (Van Iersel, Kelk [25]) *Let N be a phylogenetic network on \mathcal{X} and \mathcal{T} a set of binary trees on \mathcal{X} . Then there exists a binary phylogenetic network N' on \mathcal{X} such that (a) N' has the same reticulation number and level as N , (b) if N displays all trees in \mathcal{T} then so too does N' , (c) if N is consistent with $Tr(\mathcal{T})$ then so too is N' and (d) if N represents $Cl(\mathcal{T})$ then so too does N' .*

Suppose that we have an Algorithm A which, for each fixed $r \geq 0$, can construct in polynomial time every binary network consistent with $Tr(\mathcal{T})$ that has at most r reticulations, where \mathcal{T} is a set of binary trees on \mathcal{X} . Suppose A' is an algorithm that examines in turn every network output by A and rejects it if it does not represent $Cl(\mathcal{T})$ (such a filtering step can be done for each network in polynomial time for fixed r , see proof of Lemma 3). If there exists a network N on \mathcal{X} with at most r reticulations that represents $Cl(\mathcal{T})$ then by Lemma 8 there exists a binary network N' on \mathcal{X} with this property. Furthermore, in that case N' will by Proposition 1 be consistent with $Tr(\mathcal{T})$. The algorithm A' is thus guaranteed to eventually find N' . The practical consequence of this is that, simply by adding a filtering step, triplet software can in some cases easily be modified to work for clusters. Indeed, they can be used to determine whether a network with r reticulations that represents \mathcal{C} exists and if so to construct all binary networks with this property. As proof of concept we have taken the triplet software SIMPLISTIC (based on the ideas described in [28]), removed its biconnectedness-checking subroutine so that it generates all binary networks with up to r reticulations (and not just all simple level- $\leq r$ binary networks), and added the cluster filtering step as described above. This whole process took only one day of programming, and lead to the new software package CLUSTISTIC which implements the result described in Theorem 2. This software is available for download at <http://skelk.sdf-eu.org/clustistic>.

6 Witnesses and a natural lower bound

Now, let $\mathcal{T} = \{T_1, \dots, T_m\}$ be a set of m not necessarily binary trees on \mathcal{X} . Let S be a non-trivial ST-set of $\mathcal{Cl}(\mathcal{T})$. We know by Lemma 6 that for each $T_i \in \mathcal{T}$ there is an edge e_i or node v_i in T_i that represents S . We define a *witness* for S in T_i as follows. If T_i contains an edge $e_i = (u_i, v_i)$ that represents S , then a witness is any leaf descendant $x_i \in \mathcal{X}$ of u_i that does not appear in S , i.e. $x_i \in (\mathcal{X}(u_i) \setminus S)$. Otherwise, from Lemma 6 there exists a node v_i in T_i that represents S , and in that case a witness is any leaf descendant $x_i \in (\mathcal{X}(v_i) \setminus S)$. The only ST-sets with no witnesses are \mathcal{X} and the empty set, hence the restriction to non-trivial ST-sets. As an example consider the four trees \mathcal{T} in Figure 8. Consider the ST-set $\{1\}$ of $\mathcal{Cl}(\mathcal{T})$. In the top-left tree and bottom-left tree the only possible witness for this is taxon 5. In the top-right tree the only witness is taxon 3, and in the bottom-right tree the only witness is taxon 2.

Given a set of trees \mathcal{T} on \mathcal{X} and a non-trivial ST-set S of \mathcal{X} , let $W \subseteq \mathcal{X}$ be any subset of taxa such that, for each tree $T_i \in \mathcal{T}$, there exists $x \in W$ that is a witness for S in T_i . We call such a set a *witness set* of S in \mathcal{T} . Clearly there exist W such that $|W| \leq m$. For example, for the set of trees in Figure 8, $\{2, 3, 5\}$ is a possible witness set for $\{1\}$.

The two following simple observations are critical.

Observation 4. *Let \mathcal{T} be a set of trees on \mathcal{X} , S a non-trivial ST-set of \mathcal{X} and W a witness set of S in \mathcal{T} . Then for each $C \in \mathcal{C}(\mathcal{T})$ such that $S \subset C$, $W \cap C \neq \emptyset$.*

Proof. For each cluster $C \in \mathcal{C}(\mathcal{T})$ there is at least one edge $e = (u, v)$ in some $T_i \in \mathcal{T}$ such that e represents C . Let e_i (v_i) be the edge (node) in T_i that represents S . Given that $S \subset C$, all leaf descendants of v must also include all leaf descendants of the tail of e_i (v_i). In particular, all possible witnesses for S in T_i . \square

To understand the meaning of Observation 4 it is helpful to again consider the example of ST-set $\{1\}$ in the context of Figure 8. We see that any cluster that is a strict superset of $\{1\}$ must contain at least one of the taxa from $\{2, 3, 5\}$.

Observation 5. *Let $\mathcal{T} = \{T_1, \dots, T_m\}$ be a set of trees on \mathcal{X} such that $m \geq 2$ and let S be a non-trivial ST-set of \mathcal{C} . Let W be a smallest-cardinality witness set of S in \mathcal{T} . If $|W| = m$ then for each $C \in \mathcal{C}(\mathcal{T})$ such that $S \cap C = \emptyset$, $W \setminus C \neq \emptyset$.*

Proof. Note that it is not possible for a witness for S in a tree $T_i \in \mathcal{T}$ to also be a witness for S in $T_j \neq T_i$, because then $|W| \leq m - 1$. So each witness in W comes from a different tree in \mathcal{T} . Suppose then that there is some $C \in \mathcal{C}(\mathcal{T})$ such that $S \cap C = \emptyset$ and $W \setminus C = \emptyset$. Clearly, since $W \neq \emptyset$, $W \subseteq C$, and $|C| \geq |W| \geq 2$. Furthermore some edge e_i in some T_i represents C . Combining the fact that $S \cap C = \emptyset$ and $W \subseteq C$ leads us to the conclusion that all elements of W are possible witnesses for S in T_i . But then some element x of W is a witness for S both in T_i and also in some $T_j \neq T_i$, contradiction. \square

As mentioned in the proof of Observation 5, a witness set W (for a given ST-set S) with $m - 1$ or fewer elements exists if and only if the possible witnesses for S ranging across the different T_i are not all mutually disjoint. If a witness set with $m - 1$ or fewer elements does not exist then in the following results any witness set with m elements will turn out to be sufficient, as a consequence of Observation 5. This will become clear in due course.

6.1 A natural lower bound on $r(\mathcal{T})$ that is tight for clusters obtained from two (not necessarily binary) trees

Lemma 9. *Given a set of clusters \mathcal{C} on \mathcal{X} , there exists a maximal ST-set tree sequence (S_1, S_2, \dots, S_p) such that $p \leq r(\mathcal{C})$.*

Proof. The proof is equivalent to that of Lemma 7 but for the fact that no empty ST-set is inserted in the maximal ST-set tree sequence. \square

We define the *maximal ST-set lower bound for \mathcal{C}* (MST lower bound for short) as the cardinality of the smallest maximal ST-set tree sequence for \mathcal{C} . By Lemma 9 this is a genuine lower bound on $r(\mathcal{C})$. In general it is however a rather weak lower-bound: consider the set \mathcal{C}^i of clusters on $X^i = \{r, x_1, \dots, x_i\}$ defined by $\{\{r, x_j\} | 1 \leq j \leq i\}$. The MST lower bound for this cluster set is always 1, while $r(\mathcal{C}^i)$ rises linearly in i . However, as we shall see the tightness of the bound is to some extent correlated with the number of trees which generate the clusters, with two trees being a special case. We first require some definitions and auxiliary lemmas.

Consider a network N on \mathcal{X} . Let $\mathcal{X}' = \{\rho\} \cup \mathcal{X}$ where $\rho \notin \mathcal{X}$ is some arbitrary symbol representing the root of N . Let T be some tree on \mathcal{X}^* where $\mathcal{X}^* \cap \mathcal{X}' = \emptyset$ and let H be a subset of \mathcal{X}' . We can obtain a new network N' on $(\mathcal{X}^* \cup \mathcal{X})$ by *hanging T from H in N* . Informally N' is obtained by hanging the tree T beneath a new reticulation which has $|H|$ incoming edges, where each such incoming edge begins “just above” an element of H . Formally the transformation is as follows. First we add a new edge (r, r') to T , where r' is the root of T and r is a new node. For each $h \in H \setminus \{\rho\}$ we then subdivide the unique edge entering h ; let h_p be the new parent (with indegree and outdegree 1) of h . For each $h \in H \setminus \{\rho\}$ we then add a new edge (h_p, r) . Finally, if $\rho \in H$ we also add an edge from the root of N to r ; we call this a *root edge*. It is clear that $r(N') = r(N) + |H| - 1$.

Lemma 10. *Let $\mathcal{T} = \{T_1, \dots, T_m\}$ be a set of m trees on \mathcal{X} and let $\mathcal{C} = Cl(\mathcal{T})$. Let S be a maximum ST-set of \mathcal{C} . Let N be any network on $\mathcal{X} \setminus S$ that represents $\mathcal{C} \setminus S$. Then it is possible to extend N to obtain a new network N' that represents \mathcal{C} such that $r(N') \leq r(N) + (m - 1)$.*

Proof. Let T_S be the unique tree on taxa set S such that $\mathcal{C}(T_S) = \mathcal{C}|_S$. Let W be a minimum-cardinality witness set for S in \mathcal{T} . Clearly $1 \leq |W| \leq m$. If $|W| = m$ then we let N' be the network obtained by hanging T_S from W in N . If $|W| < m$ then we let N' be the network obtained by hanging T_S from $W \cup \{\rho\}$ in N . Clearly, $r(N') \leq r(N) + (m - 1)$. It remains only to show that N' represents \mathcal{C} . Consider any cluster $C \in \mathcal{C}$. There are three cases to consider. (1) If $C \subseteq S$ then N' clearly represents C because T_S already represented $\mathcal{C}|_S$. (2) If $S \subset C$ then consider $C' = C \setminus S$. Clearly N represents C' . By Observation 4 there exists some $w \in W \cap C$. Since $W \cap S = \emptyset$, this implies that there exists some $w \in W \cap C'$.

To see that N' represents C consider any tree displayed by N that represents C' . We can extend this tree by “switching on” the new reticulation edge that begins above w , i.e. the edge (w_p, r) , and “switching off” the remaining reticulation edges. (3) If $S \cap C = \emptyset$ then there are two subcases. (3.1) If $|W| < m$ then we can “switch on” the root edge that enters the reticulation above T_S , i.e. the edge (ρ, r) , and “switch off” all other reticulation edges entering T_S . (3.2) If $|W| = m$ then by Observation 5 there exists $w \in W \setminus C$. In N' we can thus “switch on” the new reticulation edge that begins above w , and “switch off” the rest. \square

Theorem 3. *Let $\mathcal{T} = \{T_1, \dots, T_m\}$ be a set of m trees on \mathcal{X} and let $\mathcal{C} = Cl(\mathcal{T})$. Let p be the MST lower bound for \mathcal{C} . Then $r(\mathcal{C}) \leq (m - 1)p$.*

Proof. Given a tree T and a node u of T , we denote by $\mathcal{X}(T)$ the label set of T and by T_u the subtree rooted at u . From Lemma 9 we already know that $p \leq r(\mathcal{C})$. Now, let (S_1, S_2, \dots, S_p) be a maximal ST-set tree sequence for \mathcal{C} . We will complete the proof by showing how to explicitly construct a network N with reticulation number at most $(m - 1)p$ that represents \mathcal{C} . We define \mathcal{C}_i , $1 \leq i \leq p$, as $\mathcal{C} \setminus S_1 \setminus \dots \setminus S_i$ and \mathcal{C}_0 as \mathcal{C} . By Lemma 6 it can be seen that for each i , $\mathcal{C}_i = Cl(\mathcal{T}_i)$ where \mathcal{T}_i is a set of at most m trees on $\mathcal{X} \setminus S_1 \setminus \dots \setminus S_i$ and where $\mathcal{T}_0 = \mathcal{T}$. In particular, \mathcal{T}_{i+1} can be obtained from \mathcal{T}_i as follows. Given a tree T_j in \mathcal{T}_i , let (without loss of generality) u_j be the node of T_j such that S_i is equal to the union of the clusters represented by some not necessarily strict subset of its outgoing edges. Such a u_j exists by Lemma 6. Let $Q = \{v_1, \dots, v_k\}$ be the set of children of u_j such that for each $v \in Q$, $\mathcal{X}(T_v)$ contains at least one element of S_i . The set of trees \mathcal{T}_{i+1} can be obtained from \mathcal{T}_i by computing, for each tree T_j in \mathcal{T}_i the tree $T_j \setminus T_{v_1} \dots \setminus T_{v_k}$ i.e. pruning away the subtrees corresponding to S_i and tidying up the resulting tree.

Now, consider \mathcal{C}_p . Let N_p be the unique tree such that $\mathcal{C}(N_p) = \mathcal{C}_p$; N_p will be equal to the single tree in \mathcal{T}_p . By Lemma 10 we can obtain a network N_{p-1} with $(m - 1)$ reticulations that

represents \mathcal{C}_{p-1} by taking $\mathcal{T} = \mathcal{T}_{p-1}$, $N = N_p$ and $S = S_p$ in the proof of that lemma. We iterate this process for $p-2, p-3, \dots, 1$. This lasts at most p iterations in total, and each iteration adds $(m-1)$ to the reticulation number, thus yielding a network N_0 that represents \mathcal{C} with reticulation number (at most) $(m-1)p$. \square

Corollary 9. *Let $\mathcal{T} = \{T_1, T_2\}$ be a set of two not necessarily binary trees on \mathcal{X} and let $\mathcal{C} = Cl(\mathcal{T})$. Let p be the MST lower bound for \mathcal{C} . Then $r(\mathcal{C}) = p$.*

In [25] it is shown that it is NP-hard and APX-hard to compute $r(\mathcal{C})$ where \mathcal{C} is the set of clusters obtained from two binary trees on \mathcal{X} . The following corollary is thus immediate.

Corollary 10. *The computation of the MST lower bound is NP-hard and APX-hard.*

It is interesting to note that Lemma 9 and Corollary 10 have, in some sense, already appeared in the phylogenetic network literature, albeit in the language of *recombination networks*. Specifically, in [25] we highlight that the phylogenetic network model described there (and also used here) is in a strong sense identical to the recombination network model under the assumption of an all-0 root, the infinite sites model and multiple crossover recombination. The computational lower bound described in Algorithm 3 of [18] is, taking this equivalence into account, essentially identical to the MST lower bound. In [23] it is shown that computing this bound is NP-hard, by reduction from MAX-2-SAT. The same authors also give an exponential-time dynamic programming algorithm for computing the bound because in [18] it was not explicitly indicated how this should be computed.

7 The optimality and non-optimality of Cass

The CASS algorithm for constructing simple level- k networks was presented in [27]. The algorithm was designed to produce solutions of minimum level, not of minimum reticulation number. However, when the input is a separating set \mathcal{C} of clusters on \mathcal{X} , minimizing the level or the reticulation number is equivalent. Indeed, such cluster sets have the property that any network that represents them is simple or can easily be made simple (see Lemma 1) and a simple network contains exactly one non-trivial biconnected component.

The CASS algorithm can be used as a subroutine in a divide and conquer algorithm to construct general level- k networks. We will call this more general algorithm CASS^{DC} . The basic idea of CASS^{DC} is that it transforms each connected component of $IG(\mathcal{C})$ into a tangled set of clusters, runs CASS separately on each of these tangled sets, and combines the resulting simple networks into a single final network N (Recall that tangled cluster sets are separating). The final network N has reticulation number equal to the sum of the reticulation number of the simple networks produced by CASS, and N has level equal to the maximum level ranging over all the simple networks. For more details on the divide and conquer strategy, see [27] and [12], Section 8.2.

In [27] the authors proved that if there is a level- ≤ 2 network that represents \mathcal{C} , then CASS^{DC} will find such a solution with minimal level. Here we clarify several other properties of the algorithm. On the negative side we show (using a special separating cluster set) that the CASS algorithm does not in general minimize level. On the positive side we show that when the input set \mathcal{C} is equal to $Cl(\{T_1, T_2\})$ for any two (not necessarily binary) trees, CASS^{DC} correctly minimizes level. In fact we show something even stronger: in this case CASS^{DC} also correctly constructs networks with minimum reticulation number, which in turn is exactly equal to the hybridization number of the two input trees i.e. the number of reticulations required to display the trees themselves. We conclude with several open questions regarding CASS.

7.1 Cass: the high-level idea

Let N be a network that represents a set of clusters \mathcal{C} on \mathcal{X} . Let S be a non-trivial ST-set with respect to \mathcal{C} . We say that S is *under a cut-edge* if N contains a cut-edge (u, v) such that the subnetwork rooted at v is a tree that represents $\mathcal{C}|S$.

The next two results formed (implicitly) the direct inspiration for CASS, which was designed to be a generalization of these results.

Lemma 11. *Let N be a network that represents a set of clusters \mathcal{C} on \mathcal{X} . Let S be a non-trivial ST-set with respect to \mathcal{C} . Then there exists a network N' such that $r(N') \leq r(N)$, $\ell(N') \leq \ell(N)$, S is under a cut-edge in N' and for each ST-set S' such that $S' \cap S = \emptyset$ and S' is under a cut-edge in N , S' is also under a cut-edge in N' .*

Proof. Deferred to the appendix. □

The following corollary follows from the fact that maximal ST-sets are disjoint:

Corollary 11. *Let N be a network that represents a set of clusters \mathcal{C} . There exists a network N' such that $r(N') \leq r(N)$, $\ell(N') \leq \ell(N)$ and all maximal ST-sets (with respect to \mathcal{C}) are below cut-edges.*

The pseudocode for CASS was originally given in [27]. That exposition is however rather dense and technical. See Section 8.5 of [12] for a clearer detailed description. Here we only give the core idea of the algorithm. Let us assume without loss of generality that we want to know, given a separating set of clusters, whether a simple network solution exists with reticulation number *exactly* k , for some constant k .

CASS tries to answer this by searching through the space of all maximal ST-set tree sequences of length at most k , attempting to build a network with reticulation number k from each one. It looks first at shorter maximal ST-set tree sequences, padding those of length less than k with empty ST-sets to attain a sequence of length exactly k . (As in the proof of Theorem 2, this models the situation when removing a single SBR causes the reticulation number to drop by more than 1). If there are no maximal ST-set tree sequences of length at most k then CASS will correctly report that no solutions with reticulation number k or lower exist. Hence CASS implicitly computes and incorporates the MST lower bound.

Assuming maximal ST-set tree sequences of length at most k *do* exist, CASS examines each one to determine whether it can be constructively turned into a real solution. Let $\mathcal{S} = (S_1, \dots, S_k)$ be a (possibly padded) maximal ST-set tree sequence of length- k for \mathcal{C} . As in earlier sections we define \mathcal{C}_i , $1 \leq i \leq k$, as $\mathcal{C} \setminus S_1 \setminus \dots \setminus S_i$, and we let $\mathcal{C}_0 = \mathcal{C}$. CASS does not however work with the set \mathcal{C}_i . Instead it works with \mathcal{C}'_i which is obtained from each \mathcal{C}_i by “collapsing” *every* maximal ST-set S with respect to \mathcal{C}_i into a single new “meta-taxon”. This is an *extremely* greedy step, and is in some sense an attempt to generalize Corollary 11. The informal motivation is this: we know from Corollary 11 that there exists some network N with a minimum number of reticulations such that all maximal ST-sets are under cut-edges. Suppose we guess an SBR T on $\mathcal{X}' \subset \mathcal{X}$ (where \mathcal{X}' is a maximal ST-set) of N ; we only have to make polynomially-many guesses because there are only polynomially-many maximal ST-sets. This gives a new network N' on $\mathcal{X} \setminus \mathcal{X}'$ where perhaps not all maximal ST-sets (with respect to $\mathcal{C} \setminus \mathcal{X}'$) are under cut-edges. In particular, some SBRs of N' might correspond to non-maximal ST-sets, of which there are potentially exponentially many, so how do we efficiently guess an SBR of N' ? Fortunately, we can transform N' (in the sense of Corollary 11) to obtain a new network N'' such that $r(N'') \leq r(N')$ and where all maximal ST-sets (with respect to $\mathcal{C} \setminus \mathcal{X}'$) are under cut-edges of N'' . Hence we know that we again only have to make polynomially-many guesses to locate an SBR of N'' . Furthermore, CASS assumes that these maximal ST-sets will always remain below cut-edges, so it collapses them into the aforementioned meta-taxa. We iterate this entire process k times, concluding the “inward” phase of CASS.

This assumption is important because it affects the “outward” phase of CASS, which begins immediately after completion of the inward phase. As in for example Theorem 2 the general idea is to start with a tree that represents \mathcal{C}_k and then to work backwards, first trying all pairs of edges from which to “hang back” a SBR corresponding to S_k , then all pairs of edges (of the resulting network) from which to hang back an SBR corresponding to S_{k-1} , and so on, down to S_1 . However, before hanging back each S_i it first *decollapses* the maximal ST-sets that were collapsed into meta-taxa during the corresponding iteration of the inward phase.

In Section A.2 of the appendix we will explicitly and exhaustively walk through a specific execution of the CASS algorithm and this is helpful for clarifying exactly how the algorithm works.

7.2 The Cass algorithm is not always optimal

In this section we present a counter-example proving that the CASS algorithm does not always minimize level.

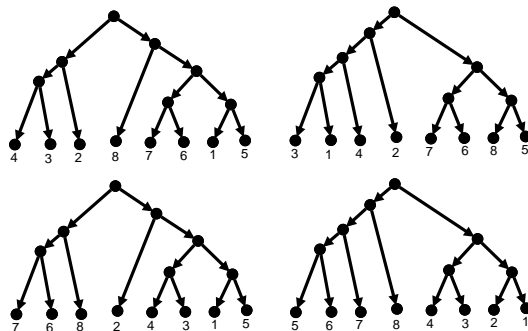


Fig. 8. Let \mathcal{T} be the set of four trees shown here. The CASS algorithm returns a network N that represents $Cl(\mathcal{T})$ where $r(N) = \ell(N) = 4$. However, Figure 9 shows that the true value of $r(Cl(\mathcal{T})) = \ell(Cl(\mathcal{T}))$ is at most 3.

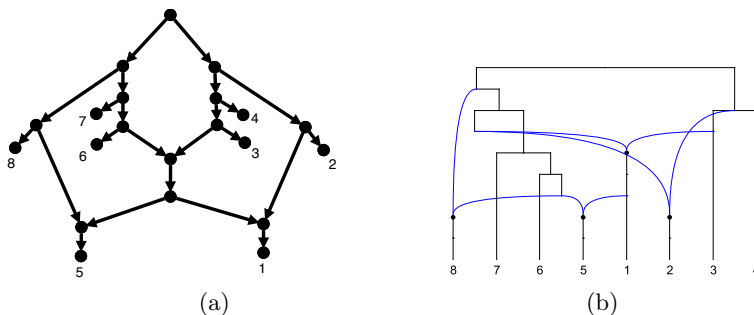


Fig. 9. (a) A simple level-3 network and (b) a simple level-4 network, both representing $Cl(\mathcal{T})$ where \mathcal{T} is defined as described in Figure 8. The level-4 network was produced by $CASS^{DC}$.

Consider the set \mathcal{T}_4 of four binary trees shown in Figure 8. It is easy to verify that the set $\mathcal{C} = Cl(\mathcal{T})$ is separating and every network that represents \mathcal{C} is thus simple or can easily be made simple. It is also easy to verify that the simple level-3 network in Figure 9(a) represents \mathcal{C} ; in the appendix we prove that this is optimal by showing that any network that represents \mathcal{C} must have reticulation number 3 or higher.

However, CASS cannot find a level-3 network. To formally prove this we show in Section A.2 of the appendix an exhaustive list of all possible executions of the algorithm with $k = 3$. CASS returns the simple level-4 network shown in Figure 9(b). To summarize, the problem with CASS seems to be that while the step of always collapsing at *every* iteration *all* maximal ST-sets and treating them as meta-taxa (in the sense of Corollary 11) is a locally optimal move, it can force us to use too many reticulation edges when hanging (trees corresponding to) maximal ST-sets back in the outward phase.

Note that this counter-example fits in a tradition of highly specific and complex counter-examples in the phylogenetic network literature. In particular we note the very similar counter-examples given initially in [8], and (based on this) in [11], which showed that one cannot minimize

reticulation number by optimizing independently over the connected components of the incompatibility graph $IG(\mathcal{C})$. The relationship between these counter-examples - all of which are linked in some way or other to simple level-3 networks - seems to be that networks with minimum reticulation number have a very subtle internal structure that seems impervious to locally optimal and greedy strategies, but that these properties only start emerging for level-3 and higher.

It is important to emphasize, however, that since we could not find *any* non-synthetic dataset for which CASS does not find an optimal solution, we still have the feeling that CASS works quite well for real data.

7.3 CASS is optimal for sets of clusters obtained from two trees

Here we show that, despite the negative news in the previous section, CASS (and more generally CASS^{DC}) correctly minimizes level in the case of clusters obtained from two not necessarily binary trees. Furthermore the algorithms also provably minimize reticulation number.

Note that the following theorem does not contradict the NP-hardness mentioned in Corollary 10 because CASS only runs in polynomial time when it bounds its search to simple level- $\leq k$ networks, for a *constant* k .

Theorem 4. *Let $\mathcal{C} = Cl(\mathcal{T})$ be a separating set of clusters where \mathcal{T} is a set of two not necessarily binary trees on \mathcal{X} . Then CASS constructs a simple network N that represents \mathcal{C} such that $\ell(N) = r(N) = \ell(\mathcal{C}) = r(\mathcal{C})$.*

Proof. Let the MST lower bound for \mathcal{C} be p . Recall that, by Corollary 9, in this case $p = r(N)$. We know that there is a maximal ST-set tree sequence (S_1, \dots, S_p) . As explained in Section 7.1 CASS will eventually find this maximal ST-set tree sequence. Now, Theorem 3 essentially works by invoking Lemma 10 p times, and in the statement of Lemma 10 there are absolutely no assumptions made about the structure of “ N ”, other than that it represents a certain set of clusters. Hence “ N ” can just as well be one of the intermediate networks constructed by CASS. It remains only to show that CASS can simulate the hanging-back construction described in the proof of Lemma 10. This is definitely so, because the proof of Lemma 10 requires the edges entering two witnesses (or the edge entering one witness and, to simulate the attachment of a root edge, the edge connecting a dummy root to the real root) to be subdivided. CASS tries subdividing all pairs of edges, including the edge between the dummy root and the real root, and hence will eventually subdivide the correct two edges. \square

The proof of Theorem 4 not only shows that CASS is optimal for sets of clusters obtained from two not necessarily binary trees, but also that in this very special case the “hanging back” (i.e. outward) phase of CASS is in some sense completely redundant. In particular: if we have already computed a maximal ST-set tree sequence, then we can easily compute a witness set for each of the maximal ST-sets in the sequence, and these witness sets directly specify a sufficient set of edges to subdivide when hanging back the maximal ST-sets. So in the case of clusters coming from two trees CASS wastes rather a lot of time trying to hang back maximal ST-sets from all possible pairs of edges, when in fact the information is already available to make this blind search unnecessary.

Theorem 4 can actually be re-formulated and extended to general (i.e. not necessarily separating) sets of clusters \mathcal{C} obtained from two not necessarily binary trees. Indeed, in Theorem 6, we will prove that, for such cluster sets, CASS^{DC} reconstructs a network N such that $r(N) = r(\mathcal{C})$ and $\ell(N) = \ell(\mathcal{C})$. However, we first need to prove Theorem 5 below, which is interesting in its own right, and several auxiliary results.

Observation 6. *Let \mathcal{C} be a set of clusters on \mathcal{X} and let $U \subseteq \mathcal{X}$ be an unseparated set with respect to \mathcal{C} . Then for any $P \subseteq \mathcal{X}$, $U \setminus P$ is unseparated with respect to $\mathcal{C}|(\mathcal{X} \setminus P)$.*

Proof. U is unseparated with respect to \mathcal{C} so for each cluster $C \in \mathcal{C}$ we have either that $C \cap U = \emptyset$, $C \subseteq U$ or $U \subseteq C$. For the case $C \cap U = \emptyset$ it is clear that $(C \setminus P) \cap (U \setminus P) = \emptyset$. For the case $C \subseteq U$ we have that $(C \setminus P) \subseteq (U \setminus P)$, and for the case $U \subseteq C$ we have that $(U \setminus P) \subseteq (C \setminus P)$. \square

Observation 7. Let \mathcal{C} be a set of clusters on \mathcal{X} and let U be the union of the set of clusters in a connected component K of $IG(\mathcal{C})$. Then U is unseparated.

Proof. Suppose U is not unseparated. Then there must exist some $C \in \mathcal{C}$ such that $C \not\subseteq U$, $U \not\subseteq C$ and $C \cap U \neq \emptyset$. Clearly C cannot be incompatible with any cluster in the connected component K , because then C would also be in the connected component K and thus $C \subseteq U$. Hence every cluster in the connected component K is either disjoint from C , or a subset of it, and there is at least one of each type of cluster because $U \setminus C \neq \emptyset$ and U is the union of all the clusters in K . However, clusters in K that are disjoint from C , are always compatible with clusters in the connected component that are contained inside C , so the connected component is not connected, contradiction. \square

Observation 8. Given a set of clusters \mathcal{C} on \mathcal{X} , let S be an ST-set with respect to \mathcal{C} and let U be an unseparated set such that $U \subseteq S$. Then U is also an ST-set for \mathcal{C} .

Proof. We only need to show that all pairs of clusters in $\mathcal{C}|U$ are compatible. Clearly, for each $C \in \mathcal{C}$ we have that $C|U \subseteq U$. Now, recall that, because U is unseparated, for each $C \in \mathcal{C}$ we have either $C \subseteq U$, $U \subseteq C$ or $C \cap U = \emptyset$. Suppose by contradiction that for some $C_1 \neq C_2 \in \mathcal{C}|U$, C_1 and C_2 are incompatible. But then $\emptyset \subset C_1, C_2 \subset U$. But in that case $C_1, C_2 \in \mathcal{C}$ and $C_1, C_2 \subset S$, contradicting the fact that S was an ST-set. \square

For a set of clusters \mathcal{C} and an unseparated set U with respect to \mathcal{C} , $\mathcal{C}_{U \rightarrow u}$ denotes the new cluster set obtained by replacing all elements of U with a single new taxon u (i.e. “collapsing” U into a single taxon).

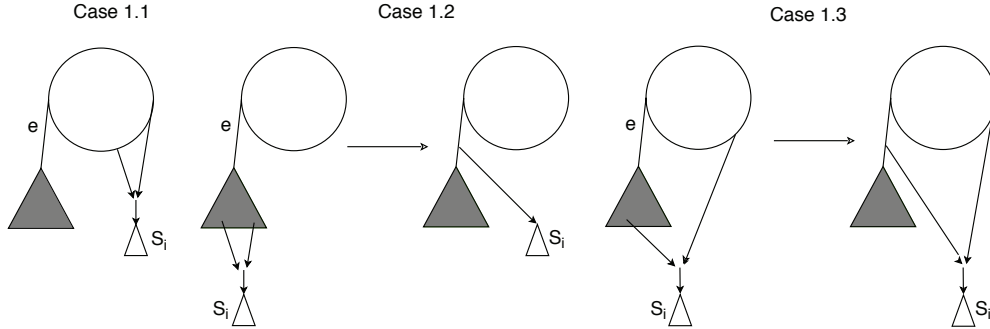


Fig. 10. The three subcases covered by Case 1, which concerns the case when the ST-set S_i that is being hung back, is disjoint from U . The parts of the network containing elements of U are depicted in grey.

Theorem 5. Let $\mathcal{T} = \{T_1, T_2\}$ be two not necessarily binary trees on \mathcal{X} , and let $\mathcal{C} = Cl(\mathcal{T})$. Let $U \subseteq \mathcal{X}$ be an unseparated set with respect to \mathcal{C} . Then $r(\mathcal{C}) = r(\mathcal{C}|U) + r(\mathcal{C}_{U \rightarrow u})$.

Proof. Let $p = r(\mathcal{C})$. We know by Corollary 9 that there exists a maximal ST-set tree sequence (S_1, \dots, S_p) . As usual we define \mathcal{C}_i , $1 \leq i \leq p$, as $\mathcal{C} \setminus S_1 \setminus \dots \setminus S_i$, and we let $\mathcal{C}_0 = \mathcal{C}$. We let $\mathcal{X}_i = \cup_{C \in \mathcal{C}_i} C$ where $\mathcal{X}_0 = \mathcal{X}$. (Note that $\mathcal{X}_i = \mathcal{X} \setminus S_1 \setminus \dots \setminus S_i$). Since $U \setminus S_1 \setminus \dots \setminus S_i = \mathcal{X}_i \cap U$, by repeated application of Observation 6, $\mathcal{X}_i \cap U$ is unseparated with respect to \mathcal{C}_i for $0 \leq i \leq p$. Now, recall (see proof of Lemma 10) that for each S_i we can identify a set of two *witnesses* (where perhaps one of the witnesses is the symbol ρ representing the root).

Here we show, for each S_i , how to hang back a tree representing $\mathcal{C}_{i-1}|S_i$ from a network N representing \mathcal{C}_i to obtain a network N' representing \mathcal{C}_{i-1} where $r(N') = r(N) + 1$ and in N' the taxa $\mathcal{X}_{i-1} \cap U$ are *exactly* the set of taxa below some cut-edge. The witnesses of S_i guide us how

to do this. If we repeat this p times we will obtain a network with p reticulations (and reticulation number p) that represents \mathcal{C} and such that the taxa in U are exactly the subset of taxa below some cut-edge. The theorem will then follow.

The first thing to do is to study the earliest point at which some elements of U are added back into the network. This is an important “base case”. Let us thus consider the largest value of i such that $\mathcal{X}_i \cap U \neq \emptyset$. Let i' be equal to this value. Now, suppose $i' = p$. We saw that by repeated application of Observation 6 $\mathcal{X}_p \cap U$ is unseparated with respect to \mathcal{C}_p . Furthermore we know that the clusters \mathcal{C}_p can be represented by a tree, so $\mathcal{X}_p \cap U$ is actually an ST-set. Hence we can assume without loss of generality (by Lemma 11) that the tree that represents \mathcal{C}_p , has a cut-edge such that $\mathcal{X}_p \cap U$ is exactly the set of taxa beneath it. Alternatively, suppose $i' < p$. In this case the first elements of U that are reintroduced into the network are a (not necessarily strict) subset of $S_{i'+1}$, so we have $\mathcal{X}_{i'} \cap U = S_{i'+1} \cap U$. Given that $\mathcal{X}_{i'} \cap U$ is unseparated w.r.t. $\mathcal{C}_{i'}$, and $S_{i'+1} \cap U$ is a subset of an ST-set, it follows by Observation 8 that $S_{i'+1} \cap U$ is also an ST-set. We may thus assume without loss of generality that $\mathcal{X}_{i'} \cap U$ is exactly the set of taxa below a cut-edge (again thanks to Lemma 11).

Henceforth we may assume that the network N that we want to hang (a tree corresponding to) $\mathcal{C}_{i-1}|S_i$ back from, contains at least one taxon of U . We will make heavy use of this fact. Let e be the cut-edge of N which the elements of $\mathcal{X}_i \cap U$ are below. Let w_1, w_2 be the two witnesses for S_i , where in some cases $w_2 = \rho$ (representing the root).

There are several cases to consider. The first case is when the tree that we are hanging back, is disjoint from U . See also Figure 10.

Case 1) $S_i \cap U = \emptyset$

Subcase 1.1) $\{w_1, w_2\} \cap U = \emptyset$. In this case we can simply hang back from $\{w_1, w_2\}$ because we are not adding any new elements of U and we are not subdividing any edge reachable from e . Hence e remains the cut-edge which all present elements of U are below.

Subcase 1.2) $\{w_1, w_2\} \subseteq U$. If we simply hang S_i back from $\{w_1, w_2\}$ then we obtain a network that represents \mathcal{C}_{i-1} . However, we see from this that every cluster C in \mathcal{C}_{i-1} that is a strict superset of S_i , must contain at least one element of U . S_i is disjoint from U so by unseparation C must also contain all other elements of U in the network. Hence we do not actually need to put S_i below a reticulation at all: we can simply attach it to a single new cut-edge that subdivides e . (This case cannot actually occur because it saves a reticulation and hence implies that $r(\mathcal{C}) < p$.)

Subcase 1.3) (wlog) $w_1 \in U, w_2 \notin U$. Suppose we simply hang back from $\{w_1, w_2\}$, obtaining a network N' that represents \mathcal{C}_{i-1} . Note that after doing this any cluster that passes through the reticulation edge starting just above w_1 , must (because of the unseparation of U) contain all elements of U that are in the network. So we can subdivide the cut-edge e and move the tail of of that reticulation edge to the newly created vertex.

The next case is when all elements of S_i are in U , see also Figure 11.

Case 2) $S_i \subseteq U$

Subcase 2.1) $\{w_1, w_2\} \cap U = \emptyset$. In this case we don't need a reticulation at all: we can just hang S_i from a single new cut-edge that subdivides e . (This case cannot actually happen because it saves a reticulation: see case 1.2).

Subcase 2.2) $\{w_1, w_2\} \subseteq U$. This case is fine because if we hang back from w_1 and w_2 all the elements of U remain below the cut-edge e .

Subcase 2.3) (wlog) $w_1 \in U, w_2 \notin U$. Suppose we hang back from w_1 and w_2 to obtain a

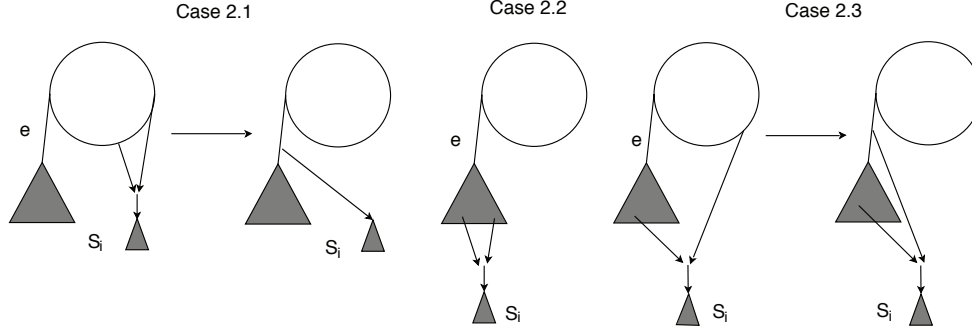


Fig. 11. The three subcases covered by Case 2, which concerns the case when the ST-set S_i that is being hung back, is a subset of U . The parts of the network containing elements of U are depicted in grey.

network that represents \mathcal{C}_{i-1} . In this case we could move the reticulation edge that starts just above w_2 (or at the root, in the case that $w_2 = \rho$) to subdivide the cut-edge e .

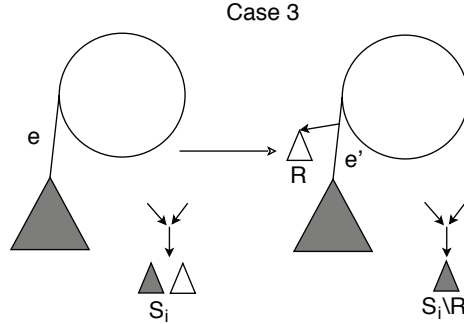


Fig. 12. Case 3, which concerns the case when the ST-set S_i that is being hung back, contains elements of U and elements not in U . The parts of the network containing elements of U are depicted in grey. Note that in this case we do not care where the reticulation edges connect to the rest of the network. Here $R = S_i \setminus U$.

The final case is where S contains at least one element of U and at least one element not in U . See also Figure 12.

Case 3) $S_i \cap U \neq \emptyset$ and $S_i \setminus U \neq \emptyset$

In this case we will apply a transformation which brings us back into Case 2. Let $R = S_i \setminus U$. Suppose we hang S_i back from its two witnesses w_1 and w_2 , to obtain a network N' that represents \mathcal{C}_{i-1} . Consider any cluster $C \in \mathcal{C}_{i-1}$ such that $C \cap R \neq \emptyset$. Then either $C \subseteq R$ or $\mathcal{X}_{i-1} \cap U \subseteq C$. (The second condition holds because, if $C \not\subseteq R$, then it must contain some element in $S_i \cap U$, and hence all elements of U in the network). So a cluster C that is a *strict* subset of S_i is either a subset of R or disjoint from R . Indeed, if C contains one element of R and one of the unseparated set U , C must contain all elements of U , a contradiction since C is a strict subset of S_i . It follows that R is unseparated so by Observation 8 R is an ST-set. Hence we may assume (without loss

of generality) that R is the taxa set of some subtree T' below a cut-edge in the tree representing $\mathcal{C}_{i-1}|S_i$ that we hung back. Now, we can prune T' and regraft it back onto the network at a new vertex obtained by subdividing e . It can be verified that after this prune/regraft move the resulting network still represents \mathcal{C}_{i-1} . It remains only to apply Case 2, taking w_1, w_2 as the witnesses and $S_i \setminus R$ as the ST-set that we want to hang back. \square

Theorem 6. *Let $\mathcal{T} = \{T_1, T_2\}$ be two not necessarily binary trees on \mathcal{X} , and let $\mathcal{C} = Cl(\mathcal{T})$. When given \mathcal{C} as input, $CASS^{DC}$ computes a level- $\ell(\mathcal{C})$ network N with reticulation number $r(\mathcal{C})$ that represents \mathcal{C} .*

Proof. Observation 7 and Theorem 5 ensure that we can analyze each connected component of the incompatibility graph $IG(\mathcal{C})$ separately, which (as mentioned) is exactly what $CASS^{DC}$ does. (To see this it is helpful to note that any subset of \mathcal{C} can also, with the possible exception of some superfluous singleton clusters, be expressed as the set of clusters in two trees).

Let K be a connected component and denote by \mathcal{C}_K the set of clusters in K . Let \mathcal{X}_K be the set of taxa equal to the union of all clusters in \mathcal{C}_K . Note that \mathcal{C}_K is not necessarily a tangled set. Indeed, while $IG(\mathcal{C}_K)$ is connected, the second property of a tangled set (every pair of taxa x and y in \mathcal{X}_K is separated by \mathcal{C}_K , Section 2) does not always hold. To ensure that the latter condition holds $CASS^{DC}$ simply computes all the maximal ST-sets $\{S_1, \dots, S_k\}$ for \mathcal{C}_K and for each of them replaces all elements of S_i with a single new taxon s_i in \mathcal{C}_K , to obtain a new cluster set \mathcal{C}'_K that is tangled.

To see that the constructed network has minimum level, observe that (by Theorem 4) $CASS$ correctly computes minimum level solutions for the \mathcal{C}'_K cluster sets mentioned above. In [27] it is proven that combining minimum-level solutions for the various \mathcal{C}'_K yields a minimum-level solution for \mathcal{C} .

We now need to prove that the constructed network has reticulation number $r(\mathcal{C})$. Since maximal ST-sets are unseparated and all mutually disjoint (see Corollary 4), it follows from Theorem 5 that:

$$r(\mathcal{C}_K) = \sum_{\substack{S \text{ is a maximal} \\ \text{ST-set of } \mathcal{C}_K}} r(\mathcal{C}_K|S) + r(\mathcal{C}'_K).$$

Since $r(\mathcal{C}_K|S)$ always equals zero when S is an ST-set, $r(\mathcal{C}_K) = r(\mathcal{C}'_K)$. Moreover, since the sets \mathcal{X}_K are also unseparated (from Observation 7) and all mutually disjoint, then from Theorem 5 we have that:

$$r(\mathcal{C}) = \sum_{\substack{K \text{ is a connected} \\ \text{component of } IG(\mathcal{C})}} r(\mathcal{C}|\mathcal{X}_K) + r(\mathcal{C}'),$$

where \mathcal{C}' is obtained from \mathcal{C} by replacing all elements of \mathcal{X}_K with a single new taxon x_K . Obviously, $r(\mathcal{C}') = 0$. Since $r(\mathcal{C}|\mathcal{X}_K) = r(\mathcal{C}_K) = r(\mathcal{C}'_K)$, and $r(\mathcal{C}'_K) = \ell(\mathcal{C}'_K)$ because \mathcal{C}'_K is separating, this concludes the proof that the constructed network has reticulation number $r(\mathcal{C})$. \square

7.4 Cass can be used to compute the hybridization number of two not necessarily binary trees

It is important to understand the relationship between the results presented in the previous section and the extensive literature on computing the hybridization distance of two *trees* [3,1,2,4,17] i.e. the problem of displaying the trees themselves, and not just their clusters.

For two binary trees $\mathcal{T} = \{T_1, T_2\}$ on \mathcal{X} the hybridization distance $h(\mathcal{T})$ is defined as the minimum reticulation number $r_t(\mathcal{T})$ of any network that displays both the trees in \mathcal{T} . In [25] we showed that $r_t(\mathcal{T}) = r(Cl(\mathcal{T}))$. As discussed in [25] that means that both positive and negative results

for the computation of $r_t(\mathcal{T})$ transfer automatically to computation of $r(Cl(\mathcal{T}))$ (for two binary trees). Negative results are NP-hardness and APX-hardness; positive results include fixed parameter tractability and running time improvements based on an increasingly deep understanding of maximum acyclic agreement forests.

Before proceeding there are some technical issues regarding the definition of hybridization number of a set \mathcal{T} of two *not necessarily binary* trees on \mathcal{X} . This can be attributed to the fact that several different definitions have appeared in the literature:

1. $h^+(\mathcal{T})$: the minimum reticulation number of N ranging over all networks N that display in a strict topological sense all the trees in \mathcal{T} . This is exactly the definition of *display* given in Section 2. In Figure 8 of [25] this strict definition was used.
2. $h^0(\mathcal{T})$: the minimum reticulation number of N ranging over all networks N that display some not necessarily binary refinement of each tree in \mathcal{T} . Recall that a (binary) *refinement* of a tree T on \mathcal{X} is any (binary) tree T' on \mathcal{X} such that $Cl(T) \subseteq Cl(T')$. (Note that a tree is generically considered to be a refinement of itself).
3. $h^-(\mathcal{T})$: the minimum reticulation number of N ranging over all networks N that display a binary refinement of each tree in \mathcal{T} . This was the definition used in [17].

Two of the definitions, h^0 and h^- , turn out to be equivalent. We clarify this, extend an equivalence result from [25] and thus show that $CASS^{DC}$ correctly computes the hybridization number of two not necessarily binary trees in the sense of h^0 , equivalently h^- . In Figure 8 of [25] a set \mathcal{T} of two trees is given, one of which is non-binary, such that $r(Cl(\mathcal{T})) < h^+(\mathcal{T})$. A result showing that $CASS$ computes h^+ was thus already excluded.

Observation 9. $h^0(\mathcal{T}) = h^-(\mathcal{T})$ for all sets \mathcal{T} of not necessarily binary trees on the same taxa set \mathcal{X} .

Proof. The fact that $h^0(\mathcal{T}) \leq h^-(\mathcal{T})$ follows immediately from the definitions. Suppose by way of contradiction that there exists a set of not necessarily binary trees \mathcal{T} such that $h^-(\mathcal{T}) > h^0(\mathcal{T})$. Let N be any network with reticulation number $h^0(\mathcal{T})$ that displays some refinement of each tree in \mathcal{T} . Now, it is not too difficult to see (using for example the transformation described in Lemma 2 of [25]) that we can create a binary network N' such that $r(N') = r(N)$ and such that N' displays a binary refinement of each of the trees in \mathcal{T} , yielding a contradiction. \square

Given a set of trees \mathcal{T} , recall that we define $r_{tr}(\mathcal{T})$ to be the minimum reticulation number of any network that displays $Tr(\mathcal{T})$ (i.e. all the rooted triplets in the input trees). We have the following result:

Theorem 7. Let $\mathcal{T} = \{T_1, T_2\}$ be two not necessarily binary trees on \mathcal{X} . Then $h^0(\mathcal{T}) = r_{tr}(\mathcal{T})$.

Proof. Obviously, for any set \mathcal{T} of not necessarily binary trees on \mathcal{X} , $r_{tr}(\mathcal{T}) \leq h^0(\mathcal{T})$. It remains to show that this inequality is always tight. Suppose that it is not always tight. Let then $\mathcal{T} = \{T_1, T_2\}$ be two smallest (in terms of the size of $|\mathcal{X}| = n$) trees such that $h^0(\mathcal{T}) > r_{tr}(\mathcal{T})$. Clearly $n > 2$. Now, let N_{tr} be any network that is consistent with $Tr(\mathcal{T}) = R$ such that $r(N_{tr}) = r_{tr}(R)$. If $r_{tr}(R) = 0$ we have a contradiction because this only occurs if T_1 and T_2 have a common refinement, in which case $h^0(\mathcal{T}) = 0$. Hence $r_{tr}(R) > 0$. This means that N_{tr} has at least one SBR with taxa set S . Now, we claim that, for each $i \in \{1, 2\}$, S is the set of taxa reachable from an edge e_i (in T_i) or as the set of taxa reachable from some subset of the children of some node u_i (in T_i). If this is not so then there exists some triplet $xy|z$ in R such that $x \notin S$ and $y, z \in S$. However, N_{tr} cannot be consistent with $xy|z$ since S is the taxa set of a SBR, which by definition sits below a cut-edge, yielding a contradiction. (If e_i exists assume without loss of generality that u_i is its head). Let $Q' = \{v_1, \dots, v_k\}$ be the set of children of u_i such that for each $v \in Q'$, $\mathcal{X}(T_v)$ contains at least one element of S . Now let T'_i be the tree $T_i \setminus T_{v_1} \dots \setminus T_{v_k}$ and let $\mathcal{T}' = \{T'_1, T'_2\}$. It is clear that $r_{tr}(\mathcal{T}') \leq r_{tr}(R) - 1$. By the assumption of minimality on $|\mathcal{X}|$ we have that $h^0(\mathcal{T}') = r_{tr}(\mathcal{T}')$. Now, let N^* be a network with a minimum reticulation number that displays some refinement of each tree in \mathcal{T}' . Observe that there exists a binary tree T_S on taxa set S such that T_S is a binary

refinement of both $T_1|S$ and $T_2|S$, otherwise $R|S$ would not be an SBR (i.e. it would contain at least one reticulation).

But we can obtain a network with $r(N^*) + 1$ reticulations that displays some refinement of T_1 and T_2 by hanging T_S back below a single new reticulation in N^* such that the tails of the two reticulation edges extend the embeddings of the refinements of T'_1 and T'_2 in N^* . Hence $h^0(\mathcal{T}) \leq r_t(\mathcal{T}') + 1 \leq r_t(\mathcal{T}) < h^0(\mathcal{T})$, contradiction. \square

The following lemma extends a result from [25]:

Lemma 12. *If \mathcal{T} consists of two not necessarily binary phylogenetic trees on the same set of taxa, $r_{tr}(\mathcal{T}) = h^0(\mathcal{T}) = r(Cl(\mathcal{T}))$.*

Proof. In [25] it is proven that for a not necessarily binary set of trees \mathcal{T} on the same set of taxa, $r_{tr}(\mathcal{T}) \leq r(Cl(\mathcal{T}))$. Now, if a network displays some refinement of a not necessarily binary tree T , then it represents all the clusters in $Cl(T)$ (and possibly more). Hence we also have that $r(Cl(\mathcal{T})) \leq h^0(\mathcal{T})$. Combining this with Theorem 7 gives the result. \square

Combining all these results we finally obtain the following theorem.

Theorem 8. *Let $\mathcal{T} = \{T_1, T_2\}$ be two not necessarily binary trees on \mathcal{X} , and let $\mathcal{C} = Cl(\mathcal{T})$. When given \mathcal{C} as input, $CASS^{DC}$ computes a network N such that $r(N) = h^0(\mathcal{T}) = h^-(\mathcal{T})$.*

8 Conclusion and open problems

The largest open problem emerging from this article is whether there exists a “reasonable” polynomial-time algorithm for constructing phylogenetic networks of bounded reticulation number (or level) that represent a given set of clusters. The result in Section 3, which shows a theoretical polynomial-time algorithm for constructing networks of bounded level, does not lend itself to a real-world implementation. On the other hand we have seen that for clusters obtained from *binary* trees a relatively simple and efficient algorithm can be used. At the moment however there is no reasonable polynomial-time algorithm for general cluster sets i.e. those obtained from sets of potentially non-binary trees. We have shown that CASS, which has a reasonable running time, is in general not optimal (although we had to explicitly engineer a highly synthetic counter-example to determine this). CASS is, however, an *extremely* greedy algorithm, in the sense that at every iteration it assumes that *all* maximal ST-sets are below cut-edges. Can we relax this assumption in some way to yield a slightly less greedy version of CASS that is optimal? A less urgent problem, but nevertheless very interesting, is the question whether CASS is optimal for clusters obtained from exactly three binary trees on \mathcal{X} .

A Appendix

A.1 Proofs deferred to the appendix

For a tree-node u with set of children Q , we say that a Q' -refinement, where $Q' \subset Q$ and $|Q'| \geq 2$, is the network obtained by deleting all edges between u and elements of Q' , adding a new node u' , adding the edge (u, u') , and adding edges between u' and each element of Q' . We say that a tree-node can be refined if there exists some Q' -refinement of it. Note that if a network N represents a set of clusters \mathcal{C} , then the network N' obtained by refining some tree-node of N still represents \mathcal{C} .

Observation 1. *Let \mathcal{C} be a separating set of clusters on \mathcal{X} . Let N be any network that represents \mathcal{C} . Then each node of N has at most one leaf child and for each cut-edge (u, v) in N , $|\mathcal{X}(v)| = 1$ or $\mathcal{X}(v) = \mathcal{X}$.*

Proof. Suppose N contains a cut-edge (u, v) such that $1 < |\mathcal{X}(v)| < |\mathcal{X}|$. Let C be any cluster represented by N . Then either $C \cap \mathcal{X}(v) = \emptyset$, $\mathcal{X}(v) \subseteq C$ or $C \subseteq \mathcal{X}(v)$, because (u, v) is a cut-edge. Hence $\mathcal{X}(v)$ is unseparated, contradiction. Now, suppose some node of N has two leaf-children $x \neq y \in \mathcal{X}$. Then every non-singleton cluster that contains x , also contains y , and vice-versa, meaning that $\{x, y\}$ is an unseparated set, contradiction. \square

Lemma 1. *Let \mathcal{C} be a separating set of clusters on \mathcal{X} . Let N be any network that represents \mathcal{C} . Then there exists a simple network N^* with at most one leaf-child per node such that $\ell(N^*) \leq \ell(N)$.*

Proof. Clearly Observation 1 holds for N . We will transform N to obtain N^* . We repeatedly and arbitrarily apply any of the following four operations until any of them can be applied. (a) If N contains a node u that has indegree and outdegree both larger than 1 then create a new node u' , add the edge (u, u') and move the tails of all edges that leave u , to u' . (b) If N contains a cut-edge (u, v) such that $\mathcal{X}(v) = \mathcal{X}$, then keep all nodes and edges reachable from v by directed paths, take v as a new root, and discard the rest of the network. (c) If N contains a cut-edge (u, v) such that $\mathcal{X}(v) = \{x\}$ where $x \in \mathcal{X}$, but v is not labelled by x , then delete all nodes and edges reachable from v by directed paths (but not v), and label v with taxon x . (d) If some tree-node of N can be refined to create a cut-edge (u, u') such that u' is not a leaf, then do that.

Note that after applying each operation the resulting network still represents \mathcal{C} . Furthermore, and less obviously, the operations do not raise the level of the network. This is clear for the “deletion” operations (b) and (c) but for operations (a) and (d) the critical observation is that two biconnected components of N overlap on at most one node, and in particular are edge disjoint. Let N^* be network obtained after no more operations (a)-(d) can be applied. Clearly, N^* represents \mathcal{C} , and Observation 1 applies to it. If N^* is simple we are done. If N^* is not simple then it must contain a cut-node u which (when removed) disconnects N^* into two or more non-trivial components, because N^* does not contain any cut-edges with this property by Observation 1. Let $P = \{p_1, \dots, p_i\}$ ($i \geq 0$) be the parent nodes of u and let $Q = \{q_1, \dots, q_j\}$ ($j \geq 1$) be the children of u . For a node $v \neq u$ let $R_u(v)$ be the set of nodes connected to v by an undirected path, after deletion of node u . We know that deleting u splits N^* into at least two non-trivial components i.e. components that contain at least one edge. Observe that each such component will be the union of one or more of the $(i + j)$ $R_u(v)$ sets obtained by taking $v \in P \cup Q$. Secondly, it is useful to note that $R_u(p_1) = R_u(p_2) = \dots = R_u(p_i)$. This follows because in a phylogenetic network every node can be reached from the root by some directed path, and hence the parents of u can still reach each other with undirected paths (via the root if necessary) after deletion of u . So at least one of the non-trivial components created by deletion of u is equal to $\cup_{v \in Q'} R_u(v)$ where $Q' \subseteq Q$. Consider the case that $Q' = Q$. Then we cannot have that $i = 0$, because this would mean that u is not a cut-node. So $i \geq 1$. If $i = 1$ then we conclude that (p_1, u) is a cut-edge in N^* and thus that u is a leaf, contradicting the assumption that $j \geq 1$. If $i \geq 2$ then $j = 1$ (because all reticulations in N^* have outdegree 1) and hence we can conclude that (u, q_1) was a cut-edge in N^* . This means that in N^* , q_1 is a leaf labelled by a single taxon and hence that $R_u(q_1)$ is a trivial component. But this contradicts the fact that at least two non-trivial components are created by deletion of u .

Consider finally the case that $Q' \subset Q$. Suppose $|Q'| = 1$ and let (wlog) q_1 be the only element of Q' . Then the edge (u, q_1) is actually a cut-edge in N^* , meaning that q_1 is a leaf labelled by a single taxon and (again) that $R_u(q_1)$ is a trivial component, contradiction. Hence $|Q'| \geq 2$ and $j \geq 3$. Note that we can construct a new network N^{**} which still represents \mathcal{C} by performing a Q' -refinement on u . Furthermore, the edge (u, u') created by the refinement must be a cut-edge. But then we could have applied a type-(d) operation to N^* , contradicting the assumption that no more applications of operations (a)-(d) were possible.

Hence we conclude that N^* is a simple network with at most one leaf-child per node that represents \mathcal{C} and, because the four operations (a)-(d) do not increase level, $\ell(N^*) \leq \ell(N)$. \square

Lemma 2. *Let N be a phylogenetic network on \mathcal{X} . Then we can transform N into a binary phylogenetic network N' such that N' has the same reticulation number and level as N and all clusters represented by N are also represented by N' .*

Proof. The following transformation is similar to one described in Lemma 2 of [25]. (a) For each reticulation node u with outdegree 2 or higher, introduce a new node u' , add an edge (u, u') and move the tails of edges that leave u , to u' . This ensures that all reticulation nodes have outdegree 1. (b) For each reticulation node u with indegree d ($d \geq 3$), we can replace u with a chain of $(d-1)$ reticulation nodes of indegree 2. For example, if u has 3 parents p_1, p_2, p_3 and child q_1 we delete u , add 2 new nodes u_1, u_2 and new edges $(p_1, u_1), (p_2, u_1), (u_1, u_2), (p_3, u_2), (u_2, q_1)$. To see that steps (a) and (b) does not increase the level of the network recall that biconnected components intersect on at most one node, and (as observed in the proof of Lemma 1) a reticulation node and all its parents are always in the same biconnected component.

(c) The only nodes we still need to consider are tree-nodes. Let u be a tree-node with outdegree d ($d \geq 3$). If at least one child of u is a leaf, but not all children of u are leaves, let q_1 be a leaf-child of u and q_2 be a non-leaf child of u . In this case we delete the edges $(u, q_1), (u, q_2)$, add a new node u' , and add new edges $(u, u'), (u', q_1), (u', q_2)$. We repeat this process until all tree-nodes u with outdegree 3 or higher are such that either all their children are leaves, or none of their children are leaves. For each tree-node u with outdegree 3 or higher whose children are all leaves, let $\mathcal{X}' = \mathcal{X}(u)$, delete all children Q of u and identify u with the root of an arbitrary binary tree on taxa set \mathcal{X}' . Now, all remaining tree-nodes with outdegree 3 or higher will be such that none of their children are leaves. The main subtlety here is to avoid a situation where, by refining a tree-node in the wrong way, we merge two biconnected components and thus raise the level of the network. To avoid this we apply the following operation as often as possible: (d) Let u be any tree-node with outdegree 3 or higher such that there exists a biconnected component K which contains u and at least two, but not all, of the edges leaving u . Let Q' be the children of u that K contains. Perform a Q' -refinement on u .

After no more applications of operation (d) are possible there is no danger of merging biconnected components by refining tree-nodes so we can apply the last refinement step: (e) we replace each remaining tree-node u with outdegree d ($d \geq 3$) by a chain of $(d-1)$ tree-nodes with outdegree 2. For example, if u has 3 children q_1, q_2, q_3 then we delete u , create two new nodes u_1, u_2 , add edges $(u_1, q_1), (u_1, u_2), (u_2, q_2), (u_2, q_3)$ and (if u had a parent) add an edge from the former parent of u to u' . This completes the transformation. \square

Lemma 4. *Let \mathcal{C} be a set of clusters on \mathcal{X} and let $S_1 \neq S_2$ be two ST-sets of \mathcal{C} . If $S_1 \cap S_2 \neq \emptyset$ then $S_1 \cup S_2$ is an ST-set.*

Proof. If $S_1 \subseteq S_2$ or $S_2 \subseteq S_1$ then we are done, so we assume neither such condition holds. Now, suppose that some cluster $C \in \mathcal{C}$ is incompatible with $S_1 \cup S_2$. We will derive a contradiction. Note that C is compatible with S_1 and S_2 since both S_1 and S_2 are ST-sets. Then we cannot have that $C \cap S_1 = \emptyset$ or $C \subseteq S_1$ because then C would be compatible with $S_1 \cup S_2$. It follows that $S_1 \subseteq C$. With the same reasoning we obtain that $S_2 \subseteq C$. So $S_1 \cup S_2 \subseteq C$. Hence C is compatible with $S_1 \cup S_2$, a contradiction. It remains to show is that any pair of clusters in $\mathcal{C}|(S_1 \cup S_2)$ are compatible. Consider any cluster C in $\mathcal{C}|(S_1 \cup S_2)$. Obviously, $C \subseteq S_1 \cup S_2$ and there are three possibilities: (i) $C = S_1 \cup S_2$, (ii) $C \subseteq S_1 \cap S_2$, or (iii) $C \cap (S_1 \cap S_2) = \emptyset$. To see this, suppose there exists a cluster $C \in \mathcal{C}|(S_1 \cup S_2)$ that violates all three possibilities. Then, without loss of generality, there is by violation of (ii) some $x \in S_1$ such that $x \notin S_2$ and $x \in C$. There is also, by violation of (iii), some $x' \in S_1 \cap S_2$ such that $x' \notin C$. This means that $S_1, S_2 \cap C \neq \emptyset$. Let C' be a cluster in \mathcal{C} such that $C'| (S_1 \cup S_2) = C$. It cannot be that $C' \subseteq S_2$ (because of x) so $S_2 \subseteq C'$ (because S_2 is an ST-set in \mathcal{C} and $S_2 \cap C \neq \emptyset$). Now, because (i) was violated it follows that there is some $x'' \in S_1$ such that $x'' \notin C$ and thus $x'' \notin C'$, so $S_1 \not\subseteq C'$. Now, since $S_2 \subseteq C'$ and we assumed that $S_2 \not\subseteq S_1$, we have that $C' \not\subseteq S_1$. But $C' \cap S_1 \neq \emptyset$, so C' is incompatible with S_1 and S_1 is not an ST-set, a contradiction. Hence (i)-(iii) are the only possibilities.

Now, suppose there are two clusters $C_1, C_2 \in \mathcal{C}|(S_1 \cup S_2)$ which are incompatible. Let C'_1 and C'_2 be the clusters in \mathcal{C} such that $C'_1| (S_1 \cup S_2) = C_1$ and $C'_2| (S_1 \cup S_2) = C_2$. Now, clearly type (i) clusters cannot cause incompatibilities. Furthermore, type (ii) and (iii) clusters are disjoint, so C_1, C_2 are both type (ii) clusters *or* are both type (iii) clusters. Suppose they are both type (ii) clusters i.e. both are entirely contained inside $S_1 \cap S_2$. If $C'_1 = C_1$ and $C'_2 = C_2$ then we obtain a

contradiction, since $S_1 \cap S_2 \subset S_1$. Specifically, we have that $C'_1, C'_2 \subset S_1$ and thus S_1 cannot be an ST-set. Hence, without loss of generality, we assume that C'_1 contains an element $x \notin S_1 \cup S_2$. But then $C'_1 \not\subseteq S_1$, $S_1 \not\subseteq C'_1$ (because $S_1 \cap S_2 \subset S_1$) and $S_1 \cap C'_1 \neq \emptyset$. Again, we conclude that S_1 is not an ST-set. The final case is that both C_1 and C_2 are type (iii) clusters. Clearly either $C_1, C_2 \subseteq S_1 \setminus S_2$ or $C_1, C_2 \subseteq S_2 \setminus S_1$, otherwise C_1 and C_2 would be compatible. Assume without loss of generality that $C_1, C_2 \subseteq S_1 \setminus S_2$. Since $S_1 \setminus S_2 \subset S_1$, if $C'_1 = C_1$ and $C'_2 = C_2$ then S_1 cannot be an ST-set, a contradiction. So assume C'_1 contains an element $x \notin S_1 \cup S_2$. So $C'_1 \not\subseteq S_1$. Furthermore, $S_1 \not\subseteq C'_1$ because $S_1 \cap S_2 \neq \emptyset$. But $S_1 \cap C'_1 \neq \emptyset$, so S_1 is not an ST-set, again a contradiction. \square

Lemma 5. *The maximal ST-sets of a set of clusters \mathcal{C} on \mathcal{X} can be computed in polynomial time..*

Proof. Start with a set \mathcal{S} of n singleton ST-sets. If there are two distinct ST-sets $S_1, S_2 \in \mathcal{S}$ such that $S_1 \cup S_2$ is an ST-set, then remove S_1 and S_2 from \mathcal{S} and add $S_1 \cup S_2$ to \mathcal{S} . Repeat this until it is no longer possible. This results in a set \mathcal{S} of ST-sets that partitions \mathcal{X} . Let \mathcal{M} be the set of maximal ST-sets of \mathcal{C} . We will show that $\mathcal{S} = \mathcal{M}$.

Observe that for each $S \in \mathcal{S}$ and each $M \in \mathcal{M}$, either $M \cap S = \emptyset$ or $S \subseteq M$. Indeed, from Lemma 4, if some S and M were incompatible then $S \cup M$ would be an ST-set too, contradicting the maximality of M . So each $M \in \mathcal{M}$ is partitioned by one or more elements from \mathcal{S} . Now, suppose some $S \in \mathcal{S}$ is not maximal. Then there exists some $M \in \mathcal{M}$ such that $S \subset M$. Furthermore, there must also exist some $S' \neq S$ such that $S' \subset M$ i.e. M is partitioned by at least two elements from \mathcal{S} . Hence we can write $M = S_1 \cup \dots \cup S_k$ where $k \geq 2$ and each $S_i \in \mathcal{S}$. We denote by \mathcal{S}' the set of ST-sets of \mathcal{S} partitioning M . Now, consider the set of clusters $\mathcal{C}|M$. All clusters in $\mathcal{C}|M$ are mutually compatible, because M is an ST-set. Furthermore, every cluster in $\mathcal{C}|M$ is compatible with every S_i contained inside M because the S_i are ST-sets and \mathcal{S}' partitions M .

Now, for a cluster $C \in \mathcal{S}$, let $nb(C)$ be the number of $S_i \in \mathcal{C}|M$ that it contains. Let $b = \min\{nb(C) | C \in \mathcal{C}|M \text{ and } nb(C) \geq 2\}$. Suppose b is not defined. Let $S_i \neq S_j$ be *any* two ST-sets of \mathcal{S}' . We claim that $S_i \cup S_j$ is an ST-set. First, we need to prove that $S_i \cup S_j$ is not separated by \mathcal{C} . Let us suppose this is not true and there exists a cluster $C \in \mathcal{C}$ such that $C \cap (S_i \cup S_j) \neq \emptyset$, $C \not\subseteq S_i \cup S_j$ and $S_i \cup S_j \not\subseteq C$. It follows that $C \not\subseteq S_i$ and $C \not\subseteq S_j$. Moreover, since $S_i \cup S_j \not\subseteq C$, C does not contain at least one among S_i and S_j , say S_i . Then $C \cap S_i = \emptyset$, otherwise S_i is not an ST-set. Then we have that $C \cap S_j \neq \emptyset$. If $S_j \not\subseteq C$ we have again a contradiction. So we must have $S_j \subset C$ and $C \cap S_i = \emptyset$. But since C cannot contain more than one ST-set of \mathcal{S}' , C separates M , a contradiction. It follows that $S_i \cup S_j$ is not separated by \mathcal{C} . Moreover, since no cluster of $\mathcal{C}|M$ contains more than one ST-set of \mathcal{S} , each cluster in $\mathcal{C}|M$ is entirely contained inside some single ST-set of \mathcal{S}' . It follows that $\mathcal{C}|(S_i \cup S_j)$ is a compatible set of clusters because in this case $\mathcal{C}|(S_i \cup S_j) = \mathcal{C}|S_i \cup \mathcal{C}|S_j$ and $S_i \cap S_j = \emptyset$. Then $S_i \cup S_j$ is an ST-set, contradicting the termination condition for the algorithm.

Now, suppose that b is well defined. Let C be any cluster in $\mathcal{C}|M$ such that $nb(C) = b$. Let $S_i \neq S_j$ be any two ST-sets in \mathcal{S}' that are contained in C . We claim that $S_i \cup S_j$ is an ST-set. We first show that $S_i \cup S_j$ is unseparated. Suppose this is not true. Then there exists $C' \in \mathcal{C}$ such that $C' \cap (S_i \cup S_j) \neq \emptyset$, $C' \not\subseteq (S_i \cup S_j)$ and $(S_i \cup S_j) \not\subseteq C'$. With the same argument as before we have (without loss of generality) that $S_j \subset C'$ and $C' \cap S_i = \emptyset$. This means that $M \not\subseteq C'$. Combined with the fact that $C' \cap M \neq \emptyset$ and that M is a maximal ST-set of \mathcal{C} , we have that $C' \subseteq M$ and thus $C' \in \mathcal{C}|M$. C' and C are compatible (because all clusters in $\mathcal{C}|M$ are mutually compatible), so combining that $C' \cap C \neq \emptyset$ and that $C \not\subseteq C'$, we have that $C' \subset C$. Observe that C' is partitioned by ST-sets from \mathcal{S}' . Now, we cannot have that $nb(C') \leq 1$ because S_j is a *strict* subset of C' . So $2 \leq nb(C') < b$, contradiction. It remains only to prove that all the clusters in $\mathcal{C}|(S_i \cup S_j)$ are mutually compatible. $S_i \cup S_j$ is unseparated by \mathcal{C} so for each $C \in \mathcal{C}$ we have that $C \cap (S_i \cup S_j) = \emptyset$, $S_i \cup S_j \subseteq C$ or $C \subseteq S_i \cup S_j$. For each $C \in \mathcal{C}$ such that $C \subseteq S_j \cup S_j$ we have either that $C = S_i \cup S_j$, $C \subseteq S_i$ or $C \subseteq S_j$ (because C is compatible with both S_i and S_j). Since $S_i \cap S_j = \emptyset$, it is not possible for $\mathcal{C}|(S_i \cup S_j)$ to contain two incompatible clusters, and again we conclude that $S_i \cup S_j$ is an ST-set. \square

Lemma 11. *Let N be a network that represents a set of clusters \mathcal{C} . Let S be a non-trivial ST-set with respect to \mathcal{C} . Then there exists a network N' such that $r(N') \leq r(N)$, $\ell(N') \leq \ell(N)$, S is under a cut-edge in N' and for each ST-set S' such that $S' \cap S = \emptyset$ and S' is under a cut-edge in N , S' is also under a cut-edge in N' .*

Proof. We obtain N' from N by the following transformation. Let x be any element of S . Recall that leaves in N always have indegree-1 and outdegree-0, so let (u, v) be the edge in N such that v is labelled by x . This is trivially a cut-edge. (a) Delete in N all taxa in S (but not the leaves they label, we will deal with this in step (c)). (b) Identify v with the root of a tree T_S on S that represents $\mathcal{C}|_S$. (c) Tidy up redundant parts of the network possibly created in step (a) by applying in an arbitrary order any of the following steps until no more can be applied: deleting any nodes with outdegree-0 that are not labelled by a taxon; suppressing any nodes with indegree-1 and outdegree-1; replacing any multi-edges with a single edge; deleting any node with indegree-0 and outdegree-1. This concludes the transformation. Let N' be the resulting network. To see that N' represents \mathcal{C} consider any $C \in \mathcal{C}$. If $C \subseteq S$ then clearly some edge in the tree we added in step (b) represents C . Otherwise there are only two possibilities (because S is unseparated): $S \cap C = \emptyset$ or $S \subset C$. In either case we know that there exists some tree T on \mathcal{X} which is displayed by N and such that T represents C . By applying steps (a)-(c) simultaneously to T we obtain a new tree T' that is displayed by N' and which represents C . The critical reason this works is that, if $C \not\subseteq S$, $x \in C$ (where x is the element we selected at the beginning of the proof) if and only if $S \subseteq C$. The tidying-up operations in step (c) clearly do not raise the reticulation number or the level of the network. Furthermore they leave untouched any taxa not in S . So any other ST-sets that were already under a cut-edge in N , are also under a cut-edge in N' . \square

A.2 Proof that Cass does not construct a simple level- ≤ 3 network for the input \mathcal{C} described in Figure 8.

Let us begin with the following observation.

Observation 10. *If a phylogenetic network N represents \mathcal{C} , then $r(N) \geq 3$.*

Proof. First, note that $r(N) \geq 2$. This follows because the clusters $\{8, 5\}, \{1, 5\}, \{5, 6\}$ are all in \mathcal{C} . A network with reticulation number r can display at most 2^r mutually non-isomorphic trees, and each tree displayed by the solution N can represent at most one of those clusters. Now, suppose some network N exists that represents \mathcal{C} such that $r(N) = 2$. All SBRs of N are single leaves, because otherwise \mathcal{C} would not be separating. Now, observe that at least one of the leaves $\{1, 5, 6, 8\}$ must be an SBR. Suppose this is not so, and let taxon $i \notin \{1, 5, 6, 8\}$ be a single leaf SBR. But then removing leaf i from N and tidying up the resulting network would give us a network N' with $r(N') \leq 1$ such that N' displays the three clusters described at the beginning of the proof; contradiction. So suppose 1 is the SBR. If we remove taxon 1 then the remaining network N' has $r(N') \leq 1$ and represents $\{8, 5\}, \{5, 6, 7\}$ and $\{3, 4, 5\}$, contradiction. If taxon 5 is the SBR then the same argument holds for $\{1, 2\}, \{3, 4, 1\}$ and $\{7, 6, 1\}$. (In fact, taxon 5 is symmetrical to taxon 1). If taxon 6 is the SBR then clusters $\{5, 8\}, \{5, 7\}, \{1, 3, 4, 5\}$ cause the contradiction. If taxon 8 is the SBR then clusters $\{1, 3, 4\}, \{1, 2\}, \{1, 5\}$ cause the contradiction. \square

To prove the overall result we thus exhaustively consider all maximal ST-set tree sequences of length exactly 3. These sequences can either be derived by hand, or generated computationally. Many maximal ST-set sequences of length 3 do not remove all incompatibilities in the clusters i.e. they are not tree sequences. We can thus automatically exclude these sequences from our analysis (because CASS will immediately conclude that it is not possible to construct a simple level- ≤ 3 network that represents \mathcal{C} from such a sequence). Additionally, many of the *tree* sequences can be excluded by using an argument similar to that used in Observation 10. In particular, if two maximal ST-sets have already been removed, but the input still contains three mutually incompatible clusters, then we discard this tree sequence. This is because any network that represents the remaining clusters (i.e. those obtained after removing the first two maximal ST-sets) must have

reticulation number at least 2, and subsequently hanging back the first two maximal ST-sets below reticulations raises the reticulation number of the network to at least 4. In other words, all attempts of CASS to construct a simple level- ≤ 3 network representing \mathcal{C} using this maximal ST-set tree sequence are doomed to fail.

It turns out that after applying these two filtering rules⁶ there only remains a very small number of cases to consider. These are :

1. $(\{1\}, \{5\}, \emptyset)$
2. $(\{8\}, \{1\}, \{5\})$
3. $(\{1\}, \{8\}, \{5\})$
4. $(\{5\}, \{1\}, \emptyset)$
5. $(\{2\}, \{5\}, \{1\})$
6. $(\{5\}, \{2\}, \{1\})$

Note that, because of symmetries in the cluster set, sequence 4 is entirely symmetrical to sequence 1, sequence 5 is entirely symmetrical to sequence 2 and sequence 6 is entirely symmetrical to sequence 3. Hence we can actually restrict our attention to only sequences 1-3.

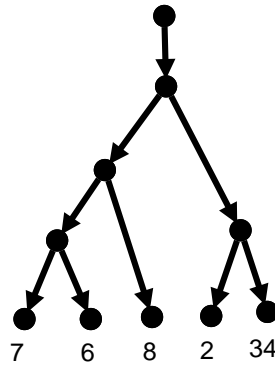


Fig. 13. The “base tree” T in the case $(\{1\}, \{5\}, \emptyset)$.

Case $(\{1\}, \{5\}, \emptyset)$

We can see already that this is a potentially important tree sequence because it is closely linked to the network shown in Figure 9(a). Specifically: remove taxon 1, causing the reticulation number to drop by 1, and then remove taxon 5 which causes the reticulation number to drop by 2 (because one of the two remaining reticulations is a child of the other). Here we show why CASS nevertheless cannot reconstruct the network.

Before removing taxon 1 the maximal ST-sets are $\{i\}$, $1 \leq i \leq 8$. After removing taxon 1 the maximal ST-sets are $\{2\}$, $\{3, 4\}$, $\{5\}$, $\{6\}$, $\{7\}$, $\{8\}$. $\{3, 4\}$ is a (new) non-singleton maximal ST-set so the CASS algorithm collapses taxa 3 and 4 into a new meta-taxon, let us call this 34 for simplicity. After collapsing the maximal ST-sets are thus $\{2\}$, $\{34\}$, $\{5\}$, $\{6\}$, $\{7\}$, $\{8\}$. Subsequently taxon 5 is removed and there are no more incompatible clusters in the input. We construct the unique tree T that represents exactly the remaining clusters (and add a dummy root): see Figure 13. For each taxon i in T let e_i be the edge in T that feeds into it.

Note how meta-taxon 34 is still collapsed. The constructive phase of CASS will now proceed as follows: (a) try and hang back a dummy taxon (i.e. the empty set \emptyset) from T below a reticulation;

⁶ The automated case-analysis can be downloaded from <http://skelk.sdf-eu.org/clustistic/caseanalysis.txt>

(b) try and hang taxon 5 back below a second reticulation; (c) decollapse meta-taxon 34; (d) try and hang taxon 1 back below a third reticulation. Now, after decollapsing meta-taxon 34 we have a network N' on taxa set $\{2, 3, 4, 5, 6, 7, 8\}$ with a cherry on taxa $\{3, 4\}$. We argue that in steps (a) and (b) certain edges of T will definitely have been subdivided by the tails of reticulation edges. Note that e_6 has definitely been subdivided to make cluster $\{5, 6\}$ possible. Similarly edge e_8 will have been subdivided to make cluster $\{5, 8\}$ possible. And edge e_{34} will have been subdivided to make cluster $\{5, 34\}$ possible. Steps (a) and (b) can subdivide at most four edges of T , but observe that if four edges are subdivided then one of the two reticulations in N' will not have been “used” i.e. there will only hang a dummy taxon beneath it. Hence N' cannot simultaneously represent $\{5, 6\}$, $\{5, 8\}$ and $\{5, 3, 4\}$, because a network with reticulation number at least 2 is required to resolve these incompatibilities. Steps (c) and (d) cannot remedy this so we can exclude the case that four edges of T have been subdivided. Hence we can safely conclude that edge e_2 of T has not yet been subdivided in N' .

Consider step (d). We have to hang back taxon 1 beneath a single reticulation i.e. a reticulation with two incoming edges. Observe that at least one of those two reticulation edges subdivides (i.e. starts from) the edge entering taxon 3, otherwise the cluster $\{1, 3\} \in \mathcal{C}$ cannot be represented. Now, note also that $\{1, 2\} \in \mathcal{C}$. To obtain this cluster the edge e_2 has to be subdivided, because otherwise every cluster that contains both 1 and 2 will also contain 3 and 4. But if both the two reticulation edges are used in this way, then any cluster that contains a taxon from $\{6, 7, 8\}$ and contains taxon 1, must also contain $\{2, 3, 4\}$. This is not satisfactory, however, because cluster $\{1, 5, 6, 7\}$ is in \mathcal{C} . In other words: at least three reticulation edges are needed to hang back taxon 1, yielding a contradiction.

It will be clear from this example that the problem with CASS is that it collapses maximal ST-set $\{3, 4\}$, forcing the iteration in which taxon 1 is hung back to use at least 3 reticulation edges instead of 2.

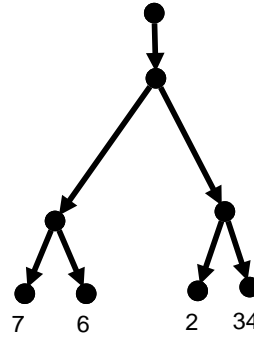


Fig. 14. The “base tree” T for the cases $(\{8\}, \{1\}, \{5\})$ and $(\{1\}, \{8\}, \{5\})$.

Case $(\{8\}, \{1\}, \{5\})$

Let T be the tree obtained after removing these three maximal ST-sets. After removing the maximal ST-set $\{1\}$ taxa 3 and 4 will, as in the previous case, have been collapsed into the meta-taxon 34. Hence T will look as in Figure 14. As in the previous case let e_i be the edge feeding into taxon i . Observe that to hang back maximal ST-set $\{5\}$ it is essential to subdivide e_6 and e_{34} ; this is the only way to ensure that the clusters $\{5, 6\}$ and $\{5, 34\}$ are represented. Let N' be the network obtained by doing this and subsequently expanding meta-taxon 34 into a cherry on taxa $\{3, 4\}$. Now, we need to hang maximal ST-set $\{1\}$ back from N' increasing the reticulation number exactly by one. Note that in any case the edge entering taxon 3 in N' will have to be subdivided (to represent the cluster $\{1, 3\}$), and also the edge entering taxon 2 will need to be subdivided (to represent the cluster $\{1, 2\}$). However, the cluster $\{1, 5, 6, 7\}$ will then definitely not

be represented. Hence we conclude that actually at least three reticulation edges are required to hang back $\{1\}$, and thus that it is impossible to hang back $\{1\}$ increasing the reticulation number exactly by one.

Case $(\{1\}, \{8\}, \{5\})$

The case has the same base tree T as the case $(\{8\}, \{1\}, \{5\})$, and again in this case there is only one way to hang back maximal ST-set $\{5\}$ as an SBR, yielding the same network N' . Now, we first need to hang back maximal ST-set $\{8\}$ in such a way that the reticulation number rises by exactly 1. Let l be the child of the real root of N' which lies on a directed path from the real root to taxon 6. We say that an edge of N' is a *left edge* if there is a directed path from the real root that contains both the tail and head of the edge and which also passes through l . Observe that when hanging back $\{8\}$ at least one left edge of N' has to be subdivided by a reticulation edge, otherwise the cluster $\{5, 6, 7, 8\}$ will not be represented by the resulting network N'' . Suppose the second reticulation edge (when hanging back $\{8\}$) subdivided neither the edge entering taxon 2, nor the edge entering taxon 3, in N' . But then, when subsequently hanging back $\{1\}$ from N'' , both these edges will have to be subdivided (to ensure that the clusters $\{1, 2\}$ and $\{1, 3\}$ are represented), and such a network cannot possibly represent the cluster $\{1, 5, 6, 7\}$. Suppose then that, when hanging back $\{8\}$, the edge entering taxon 2 in N' was subdivided by a reticulation edge e . But when hanging back $\{1\}$ the edge entering taxon 3 in N'' will definitely have to be subdivided, and also the edge entering taxon 2 or (alternatively) e . But, again, the cluster $\{1, 5, 6, 7\}$ is not represented by such a network. An essentially identical argument applies if, when hanging back $\{8\}$, only the edge entering taxon 3 was subdivided. Hence hanging back $\{8\}$ and then $\{1\}$ causes the reticulation number to rise by at least 3, rendering it impossible to construct a simple level- ≤ 3 network.

References

1. M. Bordewich, S. Linz, K. St. John, and C. Semple. A reduction algorithm for computing the hybridization number of two trees. *Evolutionary Bioinformatics*, 3:86–98, 2007.
2. M. Bordewich and C. Semple. Computing the hybridization number of two phylogenetic trees is fixed-parameter tractable. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 4(3):458–466, 2007.
3. M. Bordewich and C. Semple. Computing the minimum number of hybridization events for a consistent evolutionary history. *Discrete Applied Mathematics*, 155(8):914–928, 2007.
4. J. Collins, S. Linz, and C. Semple. Quantifying hybridization in realistic time. *Journal of Computational Biology*, 2010. To appear.
5. Philippe Gambette, Vincent Berry, and Christophe Paul. The structure of level-k phylogenetic networks. In *Proceedings of the 20th Annual Symposium on Combinatorial Pattern Matching*, CPM '09, pages 289–300, Berlin, Heidelberg, 2009. Springer-Verlag.
6. O. Gascuel, editor. *Mathematics of Evolution and Phylogeny*. Oxford University Press, Inc., 2005.
7. O. Gascuel and M. Steel, editors. *Reconstructing Evolution: New Mathematical and Computational Advances*. Oxford University Press, USA, 2007.
8. D. Gusfield, V. Bansal, V. Bafna, and Y. Song. A decomposition theory for phylogenetic networks and incompatible characters. *Journal of Computational Biology*, 14(10):1247–1272, 2007.
9. D. Gusfield, D. Hickerson, and S. Eddhu. An efficiently computed lower bound on the number of recombinations in phylogenetic networks: Theory and empirical study. *Discrete Applied Mathematics*, 155(6-7):806–830, 2007.
10. D. H. Huson and T. H. Klöpper. Beyond galled trees - decomposition and computation of galled networks. In *Research in Computational Molecular Biology (RECOMB)*, volume 4453 of *Lecture Notes in Computer Science*, pages 211–225, 2007.
11. D. H. Huson, R. Rupp, V. Berry, P. Gambette, and C. Paul. Computing galled networks from real data. *Bioinformatics*, 25(12):i85–i93, 2009.
12. Daniel H Huson, Regula Rupp, and Celine Scornavacca. *Phylogenetic Networks: Concepts, Algorithms and Applications*. Cambridge University Press, 2011. to appear.

13. D.H. Huson and C. Scornavacca. Dendroscope 3 - a program for computing and drawing rooted phylogenetic trees and networks. In preparation. Software available from: www.dendroscope.org, 2011.
14. T.N.D. Huynh, J. Jansson, N.B. Nguyen, and W.-K. Sung. Constructing a smallest refining galled phylogenetic network. In *Research in Computational Molecular Biology (RECOMB)*, volume 3500 of *Lecture Notes in Bioinformatics*, pages 265–280, 2005.
15. J. Jansson, N. B. Nguyen, and W-K. Sung. Algorithms for combining rooted triplets into a galled phylogenetic network. *SIAM Journal on Computing*, 35(5):1098–1121, 2006.
16. Jesper Jansson and Wing-Kin Sung. Inferring a level-1 phylogenetic network from a dense set of rooted triplets. *Theoretical Computer Science*, 363(1):60–68, 2006.
17. S. Linz and C. Semple. Hybridization in non-binary trees. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 6(1):30–45, 2009.
18. S. R. Myers and R. C. Griffiths. Bounds on the minimum number of recombination events in a sample history. *Genetics*, 163:375–394, 2003.
19. L. Nakhleh. *The Problem Solving Handbook for Computational Biology and Bioinformatics*, chapter Evolutionary phylogenetic networks: models and issues. Springer, 2009.
20. C. Semple. *Reconstructing Evolution - New Mathematical and Computational Advances*, chapter Hybridization Networks. Oxford University Press, 2007.
21. C. Semple and M. Steel. *Phylogenetics*. Oxford University Press, 2003.
22. T-H. To and M. Habib. Level- k phylogenetic networks are constructable from a dense triplet set in polynomial time. In *CPM09*, volume 5577 of *LNCS*, pages 275–288, 2009.
23. Bafna V. and Bansal V. Inference about recombination from haplotype data: lower bounds and recombination hotspots. *J Comput Biol.*, 13:501–21, 2006.
24. L. J. J. van Iersel, J. C. M. Keijsper, S. M. Kelk, L. Stougie, F. Hagen, and T. Boekhout. Constructing level-2 phylogenetic networks from triplets. *IEEE/ACM Transactions on Computational Biology and Bioinformatics*, 6(4):667–681, 2009.
25. L. J. J. van Iersel and S. M. Kelk. When two trees go to war. *Journal of Theoretical Biology*, 269(1):245–255, 2011.
26. L. J. J. van Iersel, S. M. Kelk, and M. Mnich. Uniqueness, intractability and exact algorithms: Reflections on level- k phylogenetic networks. *Journal of Bioinformatics and Computational Biology*, 7(2):597–623, 2009.
27. L. J. J. van Iersel, S. M. Kelk, R. Rupp, and D. H. Huson. Phylogenetic networks do not need to be complex: Using fewer reticulations to represent conflicting clusters. *Bioinformatics*, 26:i124–i131, 2010. Special issue: Proceedings of Intelligent Systems for Molecular Biology 2010 (ISMB2010), 10th-13th September 2010, Boston USA.
28. Leo van Iersel and Steven Kelk. Constructing the simplest possible phylogenetic network from triplets. *Algorithmica*, pages 1–29, 2009. 10.1007/s00453-009-9333-0.
29. Y. Wu. Close lower and upper bounds for the minimum reticulate network of multiple phylogenetic trees. *Bioinformatics*, 26:i140–i148, 2010. Special issue: Proceedings of Intelligent Systems for Molecular Biology 2010 (ISMB2010), 10th-13th September 2010, Boston USA.
30. Y. Wu and D. Gusfield. A new recombination lower bound and the minimum perfect phylogenetic forest problem. *J. Comb. Optim.*, 16(3):229–247, 2008.
31. Y. Wu and W. Jiayin. Fast computation of the exact hybridization number of two phylogenetic trees. In *Bioinformatics Research and Applications (ISBRA)*, volume 6053, pages 203–214, 2010.